

# DonorsChoose

DonorsChoose.org receives hundreds of thousands of project proposals each year for classroom projects in need of funding. Right now, a large number of volunteers is needed to manually screen each submission before it's approved to be posted on the DonorsChoose.org website.

Next year, DonorsChoose.org expects to receive close to 500,000 project proposals. As a result, there are three main problems they need to solve:

- How to scale current manual processes and resources to screen 500,000 projects so that they can be posted as quickly and as efficiently as possible
- How to increase the consistency of project vetting across different volunteers to improve the experience for teachers
- How to focus volunteer time on the applications that need the most assistance

The goal of the competition is to predict whether or not a DonorsChoose.org project proposal submitted by a teacher will be approved, using the text of project descriptions as well as additional metadata about the project, teacher, and school. DonorsChoose.org can then use this information to identify projects most likely to need further review before approval.

## About the DonorsChoose Data Set

The `train.csv` data set provided by DonorsChoose contains the following features:

Feature	Description
<code>project_id</code>	A unique identifier for the proposed project. <b>Example:</b> p036502
<code>project_title</code>	Title of the project. <b>Examples:</b> <ul style="list-style-type: none"><li>• Art Will Make You Happy!</li><li>• First Grade Fun</li></ul>
<code>project_grade_category</code>	Grade level of students for which the project is targeted. One of the following enumerated values: <ul style="list-style-type: none"><li>• Grades PreK-2</li><li>• Grades 3-5</li><li>• Grades 6-8</li><li>• Grades 9-12</li></ul>
<code>project_subject_categories</code>	One or more (comma-separated) subject categories for the project from the following enumerated list of values: <ul style="list-style-type: none"><li>• Applied Learning</li><li>• Care &amp; Hunger</li><li>• Health &amp; Sports</li><li>• History &amp; Civics</li><li>• Literacy &amp; Language</li><li>• Math &amp; Science</li><li>• Music &amp; The Arts</li><li>• Special Needs</li><li>• Warmth</li></ul> <b>Examples:</b> <ul style="list-style-type: none"><li>• Music &amp; The Arts</li><li>• Literacy &amp; Language, Math &amp; Science</li></ul>
<code>school_state</code>	State where school is located ( <a href="#">Two-letter U.S. postal code</a> ). <b>Example:</b> WY
<code>project_subject_subcategories</code>	One or more (comma-separated) subject subcategories for the project. <b>Examples:</b> <ul style="list-style-type: none"><li>• Literacy</li></ul>

Feature	Description
<code>project_resource_summary</code>	An explanation of the resources needed for the project. <b>Example:</b> <ul style="list-style-type: none"> <li>My students need hands on literacy materials to manage sensory needs!</li> </ul>
<code>project_essay_1</code>	First application essay*
<code>project_essay_2</code>	Second application essay*
<code>project_essay_3</code>	Third application essay*
<code>project_essay_4</code>	Fourth application essay*
<code>project_submitted_datetime</code>	Datetime when project application was submitted. <b>Example:</b> 2016-04-28 12:43:56.245
<code>teacher_id</code>	A unique identifier for the teacher of the proposed project. <b>Example:</b> bdf8baa8fedef6bfeec7ae4ff1c15c56
<code>teacher_prefix</code>	Teacher's title. One of the following enumerated values: <ul style="list-style-type: none"> <li>nan</li> <li>Dr.</li> <li>Mr.</li> <li>Mrs.</li> <li>Ms.</li> <li>Teacher.</li> </ul>
<code>teacher_number_of_previously_posted_projects</code>	Number of project applications previously submitted by the same teacher. <b>Example:</b> 2

\* See the section **Notes on the Essay Data** for more details about these features.

Additionally, the `resources.csv` data set provides more data about the resources required for each project. Each line in this file represents a resource required by a project:

Feature	Description
<code>id</code>	A <code>project_id</code> value from the <code>train.csv</code> file. <b>Example:</b> p036502
<code>description</code>	Description of the resource. <b>Example:</b> Tenor Saxophone Reeds, Box of 25
<code>quantity</code>	Quantity of the resource required. <b>Example:</b> 3
<code>price</code>	Price of the resource required. <b>Example:</b> 9.95

**Note:** Many projects require multiple resources. The `id` value corresponds to a `project_id` in `train.csv`, so you use it as a key to retrieve all resources needed for a project:

The data set contains the following label (the value you will attempt to predict):

Label	Description
<code>project_is_approved</code>	A binary flag indicating whether DonorsChoose approved the project. A value of 0 indicates the project was not approved, and a value of 1 indicates the project was approved.

## Notes on the Essay Data

Prior to May 17, 2016, the prompts for the essays were as follows:

- `__project_essay_1__` "Introduce us to your classroom"
- `__project_essay_2__` "Tell us more about your students"
- `__project_essay_3__` "Describe how your students will use the materials you're requesting"
- `__project_essay_3__` "Close by sharing why your project will make a difference"

Starting on May 17, 2016, the number of essays was reduced from 4 to 2, and the prompts for the first 2 essays were changed to the following:

- `__project_essay_1__` "Describe your students: What makes your students special? Specific details about their background, your neighborhood, and your school are all helpful"

your neighborhood, and your school are all helpful.

- \_\_project\_essay\_2\_\_: "About your project: How will these materials make a difference in your students' learning and improve their school lives?"

For all projects with project\_submitted\_datetime of 2016-05-17 and later, the values of project\_essay\_3 and project\_essay\_4 will be NaN.

In [1]:

```
%matplotlib inline
import warnings
warnings.filterwarnings("ignore")

import sqlite3
import pandas as pd
import numpy as np
import nltk
import string
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.feature_extraction.text import TfidfTransformer
from sklearn.feature_extraction.text import TfidfVectorizer

from sklearn.feature_extraction.text import CountVectorizer
from sklearn.metrics import confusion_matrix
from sklearn import metrics
from sklearn.metrics import roc_curve, auc
from nltk.stem.porter import PorterStemmer

import re
# Tutorial about Python regular expressions: https://pymotw.com/2/re/
import string
from nltk.corpus import stopwords
from nltk.stem import PorterStemmer
from nltk.stem.wordnet import WordNetLemmatizer

from gensim.models import Word2Vec
from gensim.models import KeyedVectors
import pickle

from tqdm import tqdm
import os

from plotly import plotly
import plotly.offline as offline
import plotly.graph_objs as go
offline.init_notebook_mode()
from collections import Counter
```

```
C:\Users\samar\Anaconda3\lib\site-packages\gensim\utils.py:1197: UserWarning: detected Windows; aliasing chunkize to chunkize_serial
  warnings.warn("detected Windows; aliasing chunkize to chunkize_serial")
```

## 1.1 Reading Data

In [2]:

```
project_data = pd.read_csv('train_data.csv')
resource_data = pd.read_csv('resources.csv')
```

In [3]:

```
print("Number of data points in train data", project_data.shape)
print('-'*50)
print("The attributes of data :", project_data.columns.values)
```

```
Number of data points in train data (109248, 17)
```

```
-----
The attributes of data : ['Unnamed: 0' 'id' 'teacher_id' 'teacher_prefix' 'school_state'
 'project_submitted_datetime' 'project_grade_category'
 'project_subject_categories' 'project_subject_subcategories']
```

```
'project_title' 'project_essay_1' 'project_essay_2' 'project_essay_3'
'project_essay_4' 'project_resource_summary'
'teacher_number_of_previously_posted_projects' 'project_is_approved']
```

In [4]:

```
print("Number of data points in train data", resource_data.shape)
print(resource_data.columns.values)
resource_data.head(2)
```

Number of data points in train data (1541272, 4)  
['id' 'description' 'quantity' 'price']

Out[4]:

	id	description	quantity	price
0	p233245	LC652 - Lakeshore Double-Space Mobile Drying Rack	1	149.00
1	p069063	Bouncy Bands for Desks (Blue support pipes)	3	14.95

## 1.2 preprocessing of project\_subject\_categories

In [5]:

```
categories = list(project_data['project_subject_categories'].values)
# remove special characters from list of strings python:
https://stackoverflow.com/a/47301924/4084039

# https://www.geeksforgeeks.org/removing-stop-words-nltk-python/
# https://stackoverflow.com/questions/23669024/how-to-strip-a-specific-word-from-a-string
# https://stackoverflow.com/questions/8270092/remove-all-whitespace-in-a-string-in-python
cat_list = []
for i in categories:
    temp = ""
    # consider we have text like this "Math & Science, Warmth, Care & Hunger"
    for j in i.split(','): # it will split it in three parts ["Math & Science", "Warmth", "Care & Hunger"]
        if 'The' in j.split(): # this will split each of the category based on space "Math & Science" => "Math", "&", "Science"
            j = j.replace('The', '') # if we have the words "The" we are going to replace it with '' (i.e. removing 'The')
            j = j.replace(' ', '') # we are replacing all the ' ' (space) with '' (empty) ex: "Math & Science" => "Math&Science"
            temp += j.strip() + " " # " abc ".strip() will return "abc", remove the trailing spaces
            temp = temp.replace('&', '_') # we are replacing the & value into
    cat_list.append(temp.strip())

project_data['clean_categories'] = cat_list
project_data.drop(['project_subject_categories'], axis=1, inplace=True)

from collections import Counter
my_counter = Counter()
for word in project_data['clean_categories'].values:
    my_counter.update(word.split())

cat_dict = dict(my_counter)
sorted_cat_dict = dict(sorted(cat_dict.items(), key=lambda kv: kv[1]))
```

## 1.3 preprocessing of project\_subject\_subcategories

In [6]:

```
sub_categories = list(project_data['project_subject_subcategories'].values)
# remove special characters from list of strings python:
https://stackoverflow.com/a/47301924/4084039

# https://www.geeksforgeeks.org/removing-stop-words-nltk-python/
# https://stackoverflow.com/questions/23669024/how-to-strip-a-specific-word-from-a-string
```

```
# https://stackoverflow.com/questions/8270092/remove-all-whitespace-in-a-string-in-python

sub_cat_list = []
for i in sub_categories:
    temp = ""
    # consider we have text like this "Math & Science, Warmth, Care & Hunger"
    for j in i.split(','): # it will split it in three parts ["Math & Science", "Warmth", "Care & Hunger"]
        if 'The' in j.split(): # this will split each of the category based on space "Math & Science" => "Math", "&", "Science"
            j = j.replace('The', '') # if we have the words "The" we are going to replace it with '' (i.e removing 'The')
            j = j.replace(' ', '') # we are replacing all the ' ' (space) with '' (empty) ex: "Math & Science" => "Math&Science"
            temp += j.strip() + " #"
        temp = temp.replace('&', '_')
    sub_cat_list.append(temp.strip())

project_data['clean_subcategories'] = sub_cat_list
project_data.drop(['project_subject_subcategories'], axis=1, inplace=True)

# count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
my_counter = Counter()
for word in project_data['clean_subcategories'].values:
    my_counter.update(word.split())

sub_cat_dict = dict(my_counter)
sorted_sub_cat_dict = dict(sorted(sub_cat_dict.items(), key=lambda kv: kv[1]))
```

## 1.4 preprocessing of project\_grade\_category

In [7]:

```
prj_grade_cat = list(project_data['project_grade_category'].values)
# remove special characters from list of strings python:
# https://stackoverflow.com/a/47301924/4084039

# https://www.geeksforgeeks.org/removing-stop-words-nltk-python/
# https://stackoverflow.com/questions/23669024/how-to-strip-a-specific-word-from-a-string
# https://stackoverflow.com/questions/8270092/remove-all-whitespace-in-a-string-in-python

prj_grade_cat_list = []
for i in prj_grade_cat:
    for j in i.split(' '): # it will split by space
        j = j.replace('Grades', '') # if we have the words "Grades" we are going to replace it with ''
        # (i.e removing 'Grades')
    prj_grade_cat_list.append(j.strip())

project_data['clean_grade'] = prj_grade_cat_list
project_data.drop(['project_grade_category'], axis=1, inplace=True)

# count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
my_counter = Counter()
for word in project_data['clean_grade'].values:
    my_counter.update(word.split())

prj_grade_cat_dict = dict(my_counter)
sorted_prj_grade_cat_dict = dict(sorted(prj_grade_cat_dict.items(), key=lambda kv: kv[1]))

project_data['clean_grade'].values
```

Out[7]:

```
array(['PreK-2', '6-8', '6-8', ..., 'PreK-2', '3-5', '6-8'], dtype=object)
```

## 1.5 preprocessing of teacher\_prefix

In [8]:

```
#tea_pfx_cat = list(project_data['teacher_prefix'].values)
tea_pfx_cat = list(project_data['teacher_prefix'].astype(str).values)
```

```
# remove special characters from list of strings python:
https://stackoverflow.com/a/47301924/4084039

# https://www.geeksforgeeks.org/removing-stop-words-nltk-python/
# https://stackoverflow.com/questions/23669024/how-to-strip-a-specific-word-from-a-string
# https://stackoverflow.com/questions/8270092/remove-all-whitespace-in-a-string-in-python

##https://stackoverflow.com/questions/52736900/how-to-solve-the-attribute-error-float-object-has-no-attribute-split-in-pyth
#vectorizer.fit(project_data['teacher_prefix'].astype(str).values)

tea_pfx_cat_list = []
for i in tea_pfx_cat:
    #for j in i.split(' '): # it will split by space
    #j=j.replace('.', '') # if we have the words "Grades" we are going to replace it with ''(i.e removing 'Grades')
    i=i.replace('.', '') # if we have the words "Grades" we are going to replace it with ''(i.e removing 'Grades')
    i=i.replace('nan', '') # if we have the words "Grades" we are going to replace it with ''(i.e removing 'Grades')
    tea_pfx_cat_list.append(i.strip())

project_data['clean_tea_pfx'] = tea_pfx_cat_list
project_data.drop(['teacher_prefix'], axis=1, inplace=True)

# count of all the words in corpus python: https://stackoverflow.com/a/22898595/4084039
my_counter = Counter()
for word in project_data['clean_tea_pfx'].values:
    my_counter.update(word.split())

tea_pfx_cat_dict = dict(my_counter)
sorted_tea_pfx_cat_dict = dict(sorted(tea_pfx_cat_dict.items(), key=lambda kv: kv[1]))

project_data['clean_tea_pfx'].values
```

Out[8]:

```
array(['Mrs', 'Mr', 'Ms', ..., 'Mrs', 'Mrs', 'Ms'], dtype=object)
```

## 1.6 Text preprocessing

In [9]:

```
# merge two column text dataframe:
project_data["essay"] = project_data["project_essay_1"].map(str) + \
    project_data["project_essay_2"].map(str) + \
    project_data["project_essay_3"].map(str) + \
    project_data["project_essay_4"].map(str)
```

In [10]:

```
project_data.head(2)
```

Out[10]:

	Unnamed: 0	id	teacher_id	school_state	project_submitted_datetime	project_title	project_title
0	160221	p253737	c90749f5d961ff158d4b4d1e7dc665fc	IN	2016-12-05 13:43:57	Educational Support for English Learners at Home	My student's English learning journey at home
1	140945	p258326	897464ce9ddc600bced1151f324dd63a	FL	2016-10-25 09:22:10	Wanted: Projector for Hungry	Our student's arrival at school

	Unnamed: 0	id	teacher_id	school_state	project_submitted_datetime	Learners project_title	lea... projec

#### 1.4.2.3 Using Pretrained Models: TFIDF weighted W2V

In [11]:

```
# printing some random reviews
print(project_data['essay'].values[0])
print("="*50)
print(project_data['essay'].values[150])
print("="*50)
print(project_data['essay'].values[1000])
print("="*50)
print(project_data['essay'].values[20000])
print("="*50)
print(project_data['essay'].values[99999])
print("="*50)
```

My students are English learners that are working on English as their second or third languages. We are a melting pot of refugees, immigrants, and native-born Americans bringing the gift of language to our school. \r\n\r\n We have over 24 languages represented in our English Learner program with students at every level of mastery. We also have over 40 countries represented with the families within our school. Each student brings a wealth of knowledge and experiences to us that open our eyes to new cultures, beliefs, and respect.\r\n\r\nThe limits of your language are the limits of your world.\r\n\r\n-Ludwig Wittgenstein Our English learner's have a strong support system at home that begs for more resources. Many times our parents are learning to read and speak English alongside of their children. Sometimes this creates barriers for parents to be able to help their child learn phonetics, letter recognition, and other reading skills.\r\n\r\n\r\nBy providing these dvd's and players, students are able to continue their mastery of the English language even if no one at home is able to assist. All families with students within the Level 1 proficiency status, will be offered to be a part of this program. These educational videos will be specially chosen by the English Learner Teacher and will be sent home regularly to watch. The videos are to help the child develop early reading skills.\r\n\r\n\r\nParents that do not have access to a dvd player will have the opportunity to check out a dvd player to use for the year. The plan is to use these videos and educational dvd's for the years to come for other EL students.\r\n\r\nnannan

=====

The 51 fifth grade students that will cycle through my classroom this year all love learning, at least most of the time. At our school, 97.3% of the students receive free or reduced price lunch. Of the 560 students, 97.3% are minority students. \r\n\r\nThe school has a vibrant community that loves to get together and celebrate. Around Halloween there is a whole school parade to show off the beautiful costumes that students wear. On Cinco de Mayo we put on a big festival with crafts made by the students, dances, and games. At the end of the year the school hosts a carnival to celebrate the hard work put in during the school year, with a dunk tank being the most popular activity. My students will use these five brightly colored Hokki stools in place of regular, stationary, 4-legged chairs. As I will only have a total of ten in the classroom and not enough for each student to have an individual one, they will be used in a variety of ways. During independent reading time they will be used as special chairs students will each use on occasion. I will utilize them in place of chairs at my small group tables during math and reading times. The rest of the day they will be used by the students who need the highest amount of movement in their life in order to stay focused on school.\r\n\r\n\r\nWhenever asked what the classroom is missing, my students always say more Hokki Stools. They can't get their fill of the 5 stools we already have. When the students are sitting in a group with me on the Hokki Stools, they are always moving, but at the same time doing their work. Anytime the students get to pick where they can sit, the Hokki Stools are the first to be taken. There are always students who head over to the kidney table to get one of the stools who are disappointed as there are not enough of them. \r\n\r\n\r\nWe ask a lot of students to sit for 7 hours a day. The Hokki stools will be a compromise that allow my students to do desk work and move at the same time. These stools will help students to meet their 60 minutes a day of movement by allowing them to activate their core muscles for balance while they sit. For many of my students, these chairs will take away the barrier that exists in schools for a child who can't sit still. nannan

=====

How do you remember your days of school? Was it in a sterile environment with plain walls, rows of desks, and a teacher in front of the room? A typical day in our room is nothing like that. I work hard to create a warm inviting themed room for my students look forward to coming to each day.\r\n\r\n\r\nMy class is made up of 28 wonderfully unique boys and girls of mixed races in Arkansas.\r\n\r\nThey attend a Title I school, which means there is a high enough percentage of free and reduced-price lunch to qualify. Our school is an "open classroom" concept, which is very unique as there are no walls separating the classrooms. These 9 and 10 year-old students are very eager learners; they are like sponges, absorbing all the information and experiences and keep on wanting more. With these resources such as the comfy red throw pillows and the whimsical nautical hanging decor and the blue fish nets, I will be able to help create the mood in our classroom setting to be one of a themed nautical environment. Creating a classroom environment is very important in the

success in each and every child's education. The nautical photo props will be used with each child as they step foot into our classroom for the first time on Meet the Teacher evening. I'll take pictures of each child with them, have them developed, and then hung in our classroom ready for their first day of 4th grade. This kind gesture will set the tone before even the first day of school! The nautical thank you cards will be used throughout the year by the students as they create thank you cards to their team groups.

Your generous donations will help me to help make our classroom a fun, inviting, learning environment from day one.

It costs lost of money out of my own pocket on resources to get our classroom ready. Please consider helping with this project to make our new school year a very successful one. Thank you!

My kindergarten students have varied disabilities ranging from speech and language delays, cognitive delays, gross/fine motor delays, to autism. They are eager beavers and always strive to work their hardest working past their limitations.

The materials we have are the ones I seek out for my students. I teach in a Title I school where most of the students receive free or reduced price lunch. Despite their disabilities and limitations, my students love coming to school and come eager to learn and explore.

Have you ever felt like you had ants in your pants and you needed to groove and move as you were in a meeting? This is how my kids feel all the time. They want to be able to move as they learn or so they say.

Wobble chairs are the answer and I love them because they develop their core, which enhances gross motor and in turn fine motor skills.

They also want to learn through games, my kids don't want to sit and do worksheets. They want to learn to count by jumping and playing. Physical engagement is the key to our success. The number toss and color and shape mats can make that happen. My students will forget they are doing work and just have the fun a 6 year old deserves.

The mediocre teacher tells. The good teacher explains. The superior teacher demonstrates. The great teacher inspires. -William A. Ward

My school has 803 students which is makeup is 97.6% African-American, making up the largest segment of the student body. A typical school in Dallas is made up of 23.2% African-American students. Most of the students are on free or reduced lunch. We aren't receiving doctors, lawyers, or engineers children from rich backgrounds or neighborhoods. As an educator I am inspiring minds of young children and we focus not only on academics but on smart, effective, efficient, and disciplined students with good character.

In our classroom we can utilize the Bluetooth for swift transitions during class. I use a speaker which doesn't amplify the sound enough to receive the message. Due to the volume of my speaker my students can't hear videos or books clearly and it isn't making the lessons as meaningful. But with the bluetooth speaker my students will be able to hear and I can stop, pause and replay it at any time.

The cart will allow me to have more room for storage of things that are needed for the day and has an extra part to it I can use. The table top chart has all of the letter, words and pictures for students to learn about different letters and it is more accessible.

In [12]:

```
# https://stackoverflow.com/a/47091490/4084039
import re

def decontracted(phrase):
    # specific
    phrase = re.sub(r"won't", "will not", phrase)
    phrase = re.sub(r"can't", "can not", phrase)

    # general
    phrase = re.sub(r"n't", " not", phrase)
    phrase = re.sub(r"'re", " are", phrase)
    phrase = re.sub(r"'s", " is", phrase)
    phrase = re.sub(r"'d", " would", phrase)
    phrase = re.sub(r"'ll", " will", phrase)
    phrase = re.sub(r"\t", " not", phrase)
    phrase = re.sub(r"'ve", " have", phrase)
    phrase = re.sub(r"'m", " am", phrase)
    return phrase
```

In [13]:

```
sent = decontracted(project_data['essay'].values[20000])
print(sent)
print("="*50)
```

My kindergarten students have varied disabilities ranging from speech and language delays, cognitive delays, gross/fine motor delays, to autism. They are eager beavers and always strive to work their hardest working past their limitations.

The materials we have are the ones I seek out for my students. I teach in a Title I school where most of the students receive free or reduced price lunch. Despite their disabilities and limitations, my students love coming to school and come eager to learn and explore.

Have you ever felt like you had ants in your pants and you needed to groove and move as you were in a meeting? This is how my kids feel all the time. They want to be able to move as they learn or so they say.

Wobble chairs are the answer and I love them because they develop their core, which enhances gross motor and in turn fine motor skills.



to move as they learn or so they say. wobble chairs are the answer and I love them because they develop their core, which enhances gross motor and in turn fine motor skills. \r\nThey also want to learn through games, my kids do not want to sit and do worksheets. They want to learn to count by jumping and playing. Physical engagement is the key to our success. The number toss and color and shape mats can make that happen. My students will forget they are doing work and just have the fun a 6 year old deserves.nannan

In [14]:

My kindergarten students have varied disabilities ranging from speech and language delays, cognitive delays, gross/fine motor delays, to autism. They are eager beavers and always strive to work their hardest working past their limitations. The materials we have are the ones I seek out for my students. I teach in a Title I school where most of the students receive free or reduced price lunch. Despite their disabilities and limitations, my students love coming to school and come eager to learn and explore. Have you ever felt like you had ants in your pants and you needed to groove and move as you were in a meeting? This is how my kids feel all the time. They want to be able to move as they learn or so they say. Wobble chairs are the answer and I love them because they develop their core, which enhances gross motor and in turn fine motor skills. They also want to learn through games, my kids do not want to sit and do worksheets. They want to learn to count by jumping and playing. Physical engagement is the key to our success. The number toss and color and shape mats can make that happen. My students will forget they are doing work and just have the fun a 6 year old deserves. nannan

In [15]:

My kindergarten students have varied disabilities ranging from speech and language delays cognitive delays gross fine motor delays to autism They are eager beavers and always strive to work their hardest working past their limitations The materials we have are the ones I seek out for my students I teach in a Title I school where most of the students receive free or reduced price lunch Despite their disabilities and limitations my students love coming to school and come eager to learn and explore Have you ever felt like you had ants in your pants and you needed to groove and move as you were in a meeting This is how my kids feel all the time The want to be able to move as they learn or so they say Wobble chairs are the answer and I love them because they develop their core which enhances gross motor and in turn fine motor skills They also want to learn through games my kids do not want to sit and do worksheets They want to learn to count by jumping and playing Physical engagement is the key to our success The number toss and color and shape mats can make that happen My students will forget they are doing work and just have the fun a 6 year old deserves nannan

In [16]:

In [17]:

In [18]:

Out [18] :

In [19]:

Out [19] :

	Unnamed: 0	id	teacher_id	school_state	project_submitted_datetime	project_title	project_description
0	160221	p253737	c90749f5d961ff158d4b4d1e7dc665fc	IN	2016-12-05 13:43:57	Educational Support for English Learners at Home	My student is struggling with English that are not in the classroom.
1	140945	p258326	807464ee0ddc600baed1151f324dd63c	FL	2016-10-25 09:33:10	Wanted: Projector for	Our student arrived in the classroom with a projector.

Unnamed: 0	id	teacher_id	school_state	project_submitted_datetime	Hungry project title Learners	school projec lea..

In [20]:

```
# printing some random essays.
print(project_data['project_title'].values[0])
print("="*50)
print(project_data['project_title'].values[150])
print("="*50)
print(project_data['project_title'].values[1000])
print("="*50)
print(project_data['project_title'].values[20000])
print("="*50)
print(project_data['project_title'].values[99999])
print("="*50)
```

```
Educational Support for English Learners at Home
=====
More Movement with Hokki Stools
=====
Sailing Into a Super 4th Grade Year
=====
We Need To Move It While We Input It!
=====
Inspiring Minds by Enhancing the Educational Experience
=====
```

In [21]:

```
sent_title = decontracted(project_data['project_title'].values[20000])
print(sent_title)
print("="*50)
```

```
We Need To Move It While We Input It!
=====
```

In [22]:

```
# \r \n \t remove from string python: http://texthandler.com/info/remove-line-breaks-python/
sent_title = sent_title.replace('\r', ' ')
sent_title = sent_title.replace('\n', ' ')
sent_title = sent_title.replace('\t', ' ')
print(sent_title)
```

```
We Need To Move It While We Input It!
```

In [23]:

```
#remove spacial character: https://stackoverflow.com/a/5843547/4084039
sent_title = re.sub('[^A-Za-z0-9]+', ' ', sent_title)
print(sent_title)
```

```
We Need To Move It While We Input It
```

In [24]:

```
# Combining all the above statemennts
from tqdm import tqdm
preprocessed_title = []
# tqdm is for printing the status bar
for sentence in tqdm(project_data['project_title'].values):
    sent_title = decontracted(sentence)
    sent_title = sent_title.replace('\r', ' ')
    sent_title = sent_title.replace('\n', ' ')
    sent_title = sent_title.replace('\t', ' ')
    sent_title = re.sub('[^A-Za-z0-9]+', ' ', sent_title)
    # https://gist.github.com/sehleier/554280
```

```
100%|██████████████████████████████████████████████████████████████████████████| 109248/109248  
[00:02<00:00, 42889.30it/s]
```

```
# after preprocessing
preprocessed_title[10]
```

'reading changes lives'

```
# Combining all the above statements
from tqdm import tqdm
preprocessed_prj_sum = []
# tqdm is for printing the status bar
for sentence in tqdm(project_data['project_resource_summary'].values):
    sent_title = decontracted(sentence)
    sent_title = sent_title.replace('\r', ' ')
    sent_title = sent_title.replace('\n', ' ')
    sent_title = sent_title.replace('\n', ' ')
    sent_title = re.sub('[^A-Za-z0-9]+', ' ', sent_title)
    # https://gist.github.com/sebleier/554280
    sent_title = ' '.join(e for e in sent_title.split() if e not in stopwords)
    preprocessed_prj_sum.append(sent_title.lower().strip())
```

```
100%|██████████████████████████████████████████████████████████████████████████████| 109248/109248  
[00:06<00:00, 18059.69it/s]
```

```
# Suggestion 5. you can try improving the score using feature engineering hacks. Try including length, summary
# and observe the results and re-submit the assignment.

# https://stackoverflow.com/questions/18827198/python-count-number-of-words-in-a-list-strings
preprocessed_essays_wc = []
for item in tqdm(preprocessed_essays):
    preprocessed_essays_wc.append(len(item.split()))

print(preprocessed_essays_wc[101])
```

```
100%|██████████████████████████████████████████████████████████████████████████| 109248/109248  
[00:01<00:00, 100095.38it/s]
```

```
# Suggestion 5. you can try improving the score using feature engineering hacks. Try including length, summary
# and observe the results and re-submit the assignment.

# https://stackoverflow.com/questions/18827198/python-count-number-of-words-in-a-list-strings
preprocessed_essays_len = []
for item in tqdm(preprocessed_essays):
    preprocessed_essays_len.append(len(item))
```

```
100%|██████████████████████████████████████████████████████████████████████████| 109248/109248  
[00:00<00:00, 1824622.70it/s]
```

### 1.8.2 Numerric feature for title

```
# Suggestion 5.you can try improving the score using feature engineering hacks.Try including length,summary
# and observe the results and re-submit the assignment.

# https://stackoverflow.com/questions/18827198/python-count-number-of-words-in-a-list-strings
preprocessed_title_wc = []
for item in tqdm(preprocessed_title):
    preprocessed_title_wc.append(len(item.split()))

print(preprocessed_title_wc[101])
```

3

```
# Suggestion 5.you can try improving the score using feature engineering hacks.Try including length,summary
# and observe the results and re-submit the assignment.

# https://stackoverflow.com/questions/18827198/python-count-number-of-words-in-a-list-strings
preprocessed_title_len = []
for item in tqdm(preprocessed_title):
    #print(preprocessed_title)
    preprocessed_title_len.append(len(item))
    #print(len(preprocessed_title))

print(preprocessed_title_len[101])
```

18

### 1.8.3 Numerric feature for project\_summary\_resource

```
# Suggestion 5. you can try improving the score using feature engineering hacks. Try including length, summary
# and observe the results and re-submit the assignment.

# https://stackoverflow.com/questions/18827198/python-count-number-of-words-in-a-list-strings
preprocessed_prj_sum_wc = []
for item in tqdm(preprocessed_prj_sum):
    preprocessed_prj_sum_wc.append(len(item.split()))

print(preprocessed_prj_sum_wc[100])
```

17

In [32]:

```
# Suggestion 5.you can try improving the score using feature engineering hacks.Try including length,summary
# and observe the results and re-submit the assignment.

# https://stackoverflow.com/questions/18827198/python-count-number-of-words-in-a-list-strings
preprocessed_prj_sum_len = []
for item in tqdm(preprocessed_prj_sum):
    preprocessed_prj_sum_len.append(len(item))

print(preprocessed_prj_sum_len[100])
```

```
100%|██████████████████████████████████████████████████████████████████████████| 109248/109248  
[00:00<00:00, 2192846.15it/s]
```

117

## 1.9 Preparing data for models

In [33]:

```
project_data.columns
```

Out [33]:

```
Index(['Unnamed: 0', 'id', 'teacher_id', 'school_state',
      'project_submitted_datetime', 'project_title', 'project_essay_1',
      'project_essay_2', 'project_essay_3', 'project_essay_4',
      'project_resource_summary',
      'teacher_number_of_previously_posted_projects', 'project_is_approved',
      'clean_categories', 'clean_subcategories', 'clean_grade',
      'clean_tea_pfx', 'essay'],
      dtype='object')
```

we are going to consider

- ```
- school_state : categorical data
- clean_categories : categorical data
- clean_subcategories : categorical data
- project_grade_category : categorical data
- teacher_prefix : categorical data

- project_title : text data
- text : text data
- project_resource_summary: text data (optinal)

- quantity : numerical (optinal)
- teacher_number_of_previously_posted_projects : numerical
- price : numerical
```

## Using Pretrained Models: Avg W2V

In [34]:

```
'''
# Reading glove vectors in python: https://stackoverflow.com/a/38230349/4084039
def loadGloveModel(gloveFile):
    print ("Loading Glove Model")
    f = open(gloveFile,'r', encoding="utf8")
    model = {}
    for line in tqdm(f):
        splitLine = line.split()
        word = splitLine[0]
        embedding = np.array([float(val) for val in splitLine[1:]])
        model[word] = embedding
```

```

    print ("Done.",len(model)," words loaded!")
    return model
model = loadGloveModel('glove.42B.300d.txt')

# =====
Output:

Loading Glove Model
1917495it [06:32, 4879.69it/s]
Done. 1917495 words loaded!

# =====

words = []
for i in preprocod_texts:
    words.extend(i.split(' '))

for i in preprocod_titles:
    words.extend(i.split(' '))
print("all the words in the coupus", len(words))
words = set(words)
print("the unique words in the coupus", len(words))

inter_words = set(model.keys()).intersection(words)
print("The number of words that are present in both glove vectors and our coupus", \
      len(inter_words), "("np.round(len(inter_words)/len(words)*100,3),"%")

words_courpus = {}
words_glove = set(model.keys())
for i in words:
    if i in words_glove:
        words_courpus[i] = model[i]
print("word 2 vec length", len(words_courpus))

# stronging variables into pickle files python: http://www.jessicayung.com/how-to-use-pickle-to-sa
ve-and-load-variables-in-python/

import pickle
with open('glove_vectors', 'wb') as f:
    pickle.dump(words_courpus, f)

'''

```

Out[34]:

```

'n# Reading glove vectors in python: https://stackoverflow.com/a/38230349/4084039\ndef
loadGloveModel(gloveFile):\n    print ("Loading Glove Model")\n    f = open(gloveFile,\r',
encoding="utf8")\n    model = {}\n    for line in tqdm(f):\n        splitLine = line.split()\n
word = splitLine[0]\n        embedding = np.array([float(val) for val in splitLine[1:]])\n        m
odel[word] = embedding\n    print ("Done.",len(model)," words loaded!")\n    return model\nmodel =
loadGloveModel(\glove.42B.300d.txt\')\n\n# =====\n\nOutput:\n    \nLoading G
love Model\n1917495it [06:32, 4879.69it/s]\nDone. 1917495 words loaded!\n\n#
=====
\n\nwords = []\nfor i in preprocod_texts:\n    words.extend(i.split(\
'))\n\nfor i in preprocod_titles:\n    words.extend(i.split(\
'))\n\nprint("all the words in the
coupus", len(words))\n\nwords = set(words)\n\nprint("the unique words in the coupus",
len(words))\n\n\ninter_words = set(model.keys()).intersection(words)\n\nprint("The number of words tha
t are present in both glove vectors and our coupus",
len(inter_words),
("np.round(len(inter_words)/len(words)*100,3),"%")\n\n\nwords_courpus = {}\n\nwords_glove =
set(model.keys())\n\nfor i in words:\n    if i in words_glove:\n        words_courpus[i] = model[i]\r
print("word 2 vec length", len(words_courpus))\n\n\n# stronging variables into pickle files python
: http://www.jessicayung.com/how-to-use-pickle-to-save-and-load-variables-in-python/\n\nimport pic
kle\n\nwith open(\glove_vectors', \wb') as f:\n    pickle.dump(words_courpus, f)\n\n\n'

```

In [35]:

```

# stronging variables into pickle files python: http://www.jessicayung.com/how-to-use-pickle-to-sa
ve-and-load-variables-in-python/
# make sure you have the glove_vectors file
with open('glove_vectors', 'rb') as f:
    model = pickle.load(f)
    glove_words = set(model.keys())

```

# Computing Sentiment Scores

In [36]:

```
## https://monkeylearn.com/sentiment-analysis/
## http://t-redactyl.io/blog/2017/04/using-vader-to-handle-sentiment-analysis-with-social-media-text.html
#
# import nltk
# from nltk.sentiment.vader import SentimentIntensityAnalyzer
#
# import nltk
# nltk.download('vader_lexicon')
#
# sid = SentimentIntensityAnalyzer()
#
# for_sentiment = 'a person is a person no matter how small dr seuss i teach the smallest students with the biggest enthusiasm \
# for learning my students learn in many different ways using all of our senses and multiple intelligences i use a wide range \
# of techniques to help all my students succeed students in my class come from a variety of different backgrounds which makes \
# for wonderful sharing of experiences and cultures including native americans our school is a caring community of successful \
# learners which can be seen through collaborative student project based learning in and out of the classroom kindergarteners \
# in my class love to work with hands on materials and have many different opportunities to practice a skill before it is \
# mastered having the social skills to work cooperatively with friends is a crucial aspect of the kindergarten curriculum \
# montana is the perfect place to learn about agriculture and nutrition my students love to role play in our pretend kitchen \
# in the early childhood classroom i have had several kids ask me can we try cooking with real food i will take their idea \
# and create common core cooking lessons where we learn important math and writing concepts while cooking delicious healthy \
# food for snack time my students will have a grounded appreciation for the work that went into making the food and knowledge \
# of where the ingredients came from as well as how it is healthy for their bodies this project would expand our learning of \
# nutrition and agricultural cooking recipes by having us peel our own apples to make homemade applesauce make our own bread \
# and mix up healthy plants from our classroom garden in the spring we will also create our own cookbooks to be printed and \
# shared with families students will gain math and literature skills as well as a life long enjoyment for healthy cooking \
# nannan'
# ss = sid.polarity_scores(for_sentiment)
#
# The end=' ' is just to say that you want a space after the end of the statement instead of a new line character.
# for k in ss:
#     print('{0}: {1}, '.format(k, ss[k]), end='')
#
# for k in ss:
#     print('{0}: {1}, '.format(k, ss[k]))
#
# we can use these 4 things as features/attributes (neg, neu, pos, compound)
# neg: 0.0, neu: 0.753, pos: 0.247, compound: 0.93
# print(type(ss))
# print(ss)
```

In [37]:

```
import nltk
from nltk.sentiment.vader import SentimentIntensityAnalyzer

import nltk
nltk.download('vader_lexicon')

sid = SentimentIntensityAnalyzer()

from tqdm import tqdm
preprocessed_sentiments = []
# tqdm is for printing the status bar
```



```
C:\Users\samar\Anaconda3\lib\site-packages\nltk\twitter\__init__.py:20: UserWarning:

The twython library has not been installed. Some functionality from the twitter package will not be available.
```

```
100%|██| 109248/109248  
[04:46<00:00, 381.81it/s]
```

```
print(type(preprocessed_sentiments))
print(preprocessed_sentiments[1:5])
# print(preprocessed_sentiments([sentiment['neg']]))
print(sentiment['neg'])
```

In [39]:

In [40]:

```
print(project_data.columns.values)
project_data['neg'].values
```

```
[ 'Unnamed: 0' 'id' 'teacher_id' 'school_state'
  'project_submitted_datetime' 'project_title' 'project_essay_1'
  'project_essay_2' 'project_essay_3' 'project_essay_4'
  'project_resource_summary' 'teacher_number_of_previously_posted_projects'
  'project_is_approved' 'clean_categories' 'clean_subcategories'
  'clean_grade' 'clean_tea_pfx' 'essay' 'neg' 'pos' 'neu' 'compound']
```

Out[40]:

```
array([0.008, 0.037, 0.058, ..., 0.    , 0.013, 0.023])
```

## Vectorizing Numerical features

In [41]:

```
price_data = resource_data.groupby('id').agg({'price': 'sum', 'quantity': 'sum'}).reset_index()
project_data = pd.merge(project_data, price_data, on='id', how='left')
```

```
print("Number of data points in train data", project_data.shape)
print('-'*50)
print("The attributes of data :", project_data.columns.values)
```

Number of data points in train data (109248, 24)

The attributes of data : ['Unnamed: 0' 'id' 'teacher\_id' 'school\_state' 'project\_submitted\_datetime' 'project\_title' 'project\_essay\_1' 'project\_essay\_2' 'project\_essay\_3' 'project\_essay\_4' 'project\_resource\_summary' 'teacher\_number\_of\_previously\_posted\_projects']

```
'project_resource_summary' 'teacher_number_of_previously_posted_projects'
'project_is_approved' 'clean_categories' 'clean_subcategories'
'clean_grade' 'clean_tea_pfx' 'essay' 'neg' 'pos' 'neu' 'compound'
'price' 'quantity']
```

## Adding word count and length column

In [42]:

```
project_data['essay_wc'] = preprocessed_essays_wc
project_data['essay_len'] = preprocessed_essays_len

project_data['title_wc'] = preprocessed_title_wc
project_data['title_len'] = preprocessed_title_len

project_data['prj_res_sum_wc'] = preprocessed_prj_sum_wc
project_data['prj_res_sum_len'] = preprocessed_prj_sum_len

print("Number of data points in train data", project_data.shape)
print('-'*50)
print("The attributes of data :", project_data.columns.values)
```

Number of data points in train data (109248, 30)

-----

The attributes of data : ['Unnamed: 0' 'id' 'teacher\_id' 'school\_state'
'project\_submitted\_datetime' 'project\_title' 'project\_essay\_1'
'project\_essay\_2' 'project\_essay\_3' 'project\_essay\_4'
'project\_resource\_summary' 'teacher\_number\_of\_previously\_posted\_projects'
'project\_is\_approved' 'clean\_categories' 'clean\_subcategories'
'clean\_grade' 'clean\_tea\_pfx' 'essay' 'neg' 'pos' 'neu' 'compound'
'price' 'quantity' 'essay\_wc' 'essay\_len' 'title\_wc' 'title\_len'
'prj\_res\_sum\_wc' 'prj\_res\_sum\_len']

## Assignment 8: DT

### 1. Apply Decision Tree Classifier(DecisionTreeClassifier) on these feature sets

- **Set 1:** categorical, numerical features + project\_title(BOW) + preprocessed\_eassay (BOW)
- **Set 2:** categorical, numerical features + project\_title(TFIDF)+ preprocessed\_eassay (TFIDF)
- **Set 3:** categorical, numerical features + project\_title(AVG W2V)+ preprocessed\_eassay (AVG W2V)
- **Set 4:** categorical, numerical features + project\_title(TFIDF W2V)+ preprocessed\_eassay (TFIDF W2V)



### 2. Hyper paramter tuning (best `depth` in range [1, 5, 10, 50, 100, 500, 100], and the best `min\_samples\_split` in range [5, 10, 100, 500])

- Find the best hyper parameter which will give the maximum [AUC](#) value
- Find the best hyper paramter using k-fold cross validation or simple cross validation data
- Use gridsearch cv or randomsearch cv or you can also write your own for loops to do this task of hyperparameter tuning

### 3. Graphviz

- Visualize your decision tree with Graphviz. It helps you to understand how a decision is being made, given a new vector.
- Since feature names are not obtained from word2vec related models, visualize only BOW & TFIDF decision trees using Graphviz
- Make sure to print the words in each node of the decision tree instead of printing its index.
- Just for visualization purpose, limit max\_depth to 2 or 3 and either embed the generated images of graphviz in your notebook, or directly upload them as .png files.

### 4. Representation of results

- You need to plot the performance of model both on train data and cross validation data for each hyper parameter, like shown in the figure 
- Once after you found the best hyper parameter, you need to train your model with it, and find the AUC on test data and plot the ROC curve on both train and test. 

- Along with plotting ROC curve, you need to print the [confusion matrix](#) with predicted and original labels of test data points
- Once after you plot the confusion matrix with the test data, get all the `false positive data points`
  - Plot the WordCloud [WordCloud](#)
  - Plot the box plot with the `price` of these `false positive data points`
  - Plot the pdf with the `teacher\_number\_of\_previously\_posted\_projects` of these `false positive data points`

## 5. [Task-2]

- Select 5k best features from features of **Set 2** using [`feature importances`](#), discard all the other remaining features and then apply any of the model of your choice i.e. (Decision tree, Logistic Regression, Linear SVM), you need to do hyperparameter tuning corresponding to the model you selected and procedure in step 2 and step 3

## 6. Conclusion

- You need to summarize the results at the end of the notebook, summarize it in the table format. To print out a table please refer to this [prettytable library link](#)

# 2. Decision Tree

In [43]:

```
##taking 50K datapoint
project_data50K=project_data[:50000]
#project_data100K=project_data[:100000]
#X=project_data100K
X=project_data50K
print(project_data50K.shape)
#print(project_data100K.shape)
print(X.shape)
```

```
(50000, 30)
(50000, 30)
```

In [44]:

```
# makes Xi as 19 column matrix, where we create the modle and Yi as single column matrix as a class label.
#y = project_data50K['project_is_approved'].values
#project_data50K.drop(['project_is_approved'], axis=1, inplace=True)
#print(y.shape)
#project_data50K.head(1)

y = project_data['project_is_approved'].values
project_data.drop(['project_is_approved'], axis=1, inplace=True)
#print(y.shape)
project_data.head(1)

#y100K=y[:100000]
#y=y100K
y50K=y[:50000]
y=y50K

#y = project_data['project_is_approved'].values
#project_data.drop(['project_is_approved'], axis=1, inplace=True)
print(y.shape)
#project_data.head(1)
```

```
(50000,)
```

In [45]:

```
#X = project_data50K
print(X.shape)
print(y.shape)
#X1K = project_data1K
#print(X1K.shape)
```

```
(50000, 30)
(50000,)
```

## 2.1 Splitting data into Train and cross validation(or test): Stratified Sampling

In [46]:

```
# please write all the code with proper documentation, and proper titles for each subsection
# go through documentations and blogs before you start coding
# first figure out what to do, and then think about how to do.
# reading and understanding error messages will be very much helpfull in debugging your code
# when you plot any graph make sure you use
# a. Title, that describes your plot, this will be very helpful to the reader
# b. Legends if needed
# c. X-axis label
# d. Y-axis label
```

In [47]:

```
# train test split | https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.train_test_split.html
# splitting Xq and Yq in Train(further into Train and CV) and Test matrix
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.33, stratify=y)
#X_train, X_cv, y_train, y_cv = train_test_split(X_train, y_train, test_size=0.33,
stratify=y_train)

print(X_train.shape, y_train.shape)
#print(X_cv.shape, y_cv.shape)
print(X_test.shape, y_test.shape)

print("="*100)
```

```
(33500, 30) (33500,)
(16500, 30) (16500,)
```

### 2.1.1 Make Data Model Ready: encoding school\_state categorical data

In [48]:

```
from sklearn.feature_extraction.text import CountVectorizer
vectorizer = CountVectorizer(min_df=10,ngram_range=(1,2), max_features=5000)
vectorizer.fit(X_train['school_state'].values) # fit has to happen only on train data

# we use the fitted CountVectorizer to convert the text to vector
X_train_state_ohe = vectorizer.transform(X_train['school_state'].values)
X_cv_state_ohe = vectorizer.transform(X_cv['school_state'].values)
X_test_state_ohe = vectorizer.transform(X_test['school_state'].values)

print("school_state After vectorizations")
print(X_train_state_ohe.shape, y_train.shape)
#print(X_cv_state_ohe.shape, y_cv.shape)
print(X_test_state_ohe.shape, y_test.shape)
print(vectorizer.get_feature_names())
print("="*100)
aa=vectorizer.get_feature_names()
```

```
school_state After vectorizations
(33500, 51) (33500,)
(16500, 51) (16500,)
['ak', 'al', 'ar', 'az', 'ca', 'co', 'ct', 'dc', 'de', 'fl', 'ga', 'hi', 'ia', 'id', 'il', 'in', 'k
s', 'ky', 'la', 'ma', 'md', 'me', 'mi', 'mn', 'mo', 'ms', 'mt', 'nc', 'nd', 'ne', 'nh', 'nj', 'nm',
'nv', 'ny', 'oh', 'ok', 'or', 'pa', 'ri', 'sc', 'sd', 'tn', 'tx', 'ut', 'va', 'vt', 'wa', 'wi', 'wv
', 'wy']
```

## 2.1.2 Make Data Model Ready: encoding clean\_categories

In [49]:

```
from sklearn.feature_extraction.text import CountVectorizer
#vectorizer = CountVectorizer(min_df=10,ngram_range=(1,2), max_features=5000)
vectorizer = CountVectorizer(vocabulary =list(sorted_cat_dict.keys()),lowercase =False,binary=True
)
vectorizer.fit(X_train['clean_categories'].values) # fit has to happen only on train data

# we use the fitted CountVectorizer to convert the text to vector
X_train_clean_ohe = vectorizer.transform(X_train['clean_categories'].values)
#X_cv_clean_ohe = vectorizer.transform(X_cv['clean_categories'].values)
X_test_clean_ohe = vectorizer.transform(X_test['clean_categories'].values)

print("clean_categories After vectorizations")
print(X_train_clean_ohe.shape, y_train.shape)
#print(X_cv_clean_ohe.shape, y_cv.shape)
print(X_test_clean_ohe.shape, y_test.shape)
print(vectorizer.get_feature_names())
print("="*100)
b=vectorizer.get_feature_names()
```

```
clean_categories After vectorizations
(33500, 9) (33500,)
(16500, 9) (16500,)
['Warmth', 'Care_Hunger', 'History_Civics', 'Music_Arts', 'AppliedLearning', 'SpecialNeeds',
'Health_Sports', 'Math_Science', 'Literacy_Language']
=====
```



## 2.1.3 Make Data Model Ready: encoding clean\_subcategories

In [50]:

```
from sklearn.feature_extraction.text import CountVectorizer
vectorizer = CountVectorizer(vocabulary =list(sorted_sub_cat_dict.keys()),lowercase =False,binary=
True)
vectorizer.fit(X_train['clean_subcategories'].values) # fit has to happen only on train data

# we use the fitted CountVectorizer to convert the text to vector
X_train_cleanSub_ohe = vectorizer.transform(X_train['clean_subcategories'].values)
#X_cv_cleanSub_ohe = vectorizer.transform(X_cv['clean_subcategories'].values)
X_test_cleanSub_ohe = vectorizer.transform(X_test['clean_subcategories'].values)

print("clean_subcategories After vectorizations")
print(X_train_cleanSub_ohe.shape, y_train.shape)
#print(X_cv_cleanSub_ohe.shape, y_cv.shape)
print(X_test_cleanSub_ohe.shape, y_test.shape)
#print(vectorizer.get_feature_names())
print("="*100)
c=vectorizer.get_feature_names()
```

```
clean_subcategories After vectorizations
(33500, 30) (33500,)
(16500, 30) (16500,)
=====
```



## 2.1.4 Make Data Model Ready: encoding project\_grade\_category

In [51]:

```
from sklearn.feature_extraction.text import CountVectorizer
vectorizer = CountVectorizer(vocabulary =list(sorted_prj_grade_cat_dict.keys()),lowercase =False,b
inary=True)
vectorizer.fit(X_train['clean_grade'].values) # fit has to happen only on train data
```

```
# we use the fitted CountVectorizer to convert the text to vector
X_train_grade_ohe = vectorizer.transform(X_train['clean_grade'].values)
#X_cv_grade_ohe = vectorizer.transform(X_cv['clean_grade'].values)
X_test_grade_ohe = vectorizer.transform(X_test['clean_grade'].values)

print("project_grade_category After vectorizations")
print(X_train_grade_ohe.shape, y_train.shape)
#print(X_cv_grade_ohe.shape, y_cv.shape)
print(X_test_grade_ohe.shape, y_test.shape)
print(vectorizer.get_feature_names())
print("="*100)
d=vectorizer.get_feature_names()
```

```
project_grade_category After vectorizations
(33500, 4) (33500,)
(16500, 4) (16500,)
['9-12', '6-8', '3-5', 'PreK-2']
=====
```

## 2.1.5 Make Data Model Ready: encoding teacher\_prefix

In [52]:

```
from sklearn.feature_extraction.text import CountVectorizer
vectorizer = CountVectorizer(vocabulary =list(sorted_tea_pfx_cat_dict.keys()), lowercase =False, binary=True)
#https://stackoverflow.com/questions/52736900/how-to-solve-the-attribute-error-float-object-has-no-attribute-split-in-pyth
vectorizer.fit(X_train['clean_tea_pfx'].astype(str).values) # fit has to happen only on train data

# we use the fitted CountVectorizer to convert the text to vector
X_train_teacher_ohe = vectorizer.transform(X_train['clean_tea_pfx'].astype(str).values)
#X_cv_teacher_ohe = vectorizer.transform(X_cv['clean_tea_pfx'].astype(str).values)
X_test_teacher_ohe = vectorizer.transform(X_test['clean_tea_pfx'].astype(str).values)

print("teacher_prefix After vectorizations")
print(X_train_teacher_ohe.shape, y_train.shape)
#print(X_cv_teacher_ohe.shape, y_cv.shape)
print(X_test_teacher_ohe.shape, y_test.shape)
print(vectorizer.get_feature_names())
print("="*100)
e=vectorizer.get_feature_names()
```

```
teacher_prefix After vectorizations
(33500, 5) (33500,)
(16500, 5) (16500,)
['Dr', 'Teacher', 'Mr', 'Ms', 'Mrs']
=====
```

## 2.1.6 Make Data Model Ready: encoding project\_resource\_summary

In [53]:

```
vectorizer = CountVectorizer(min_df=10,ngram_range=(1,2))
vectorizer.fit(X_train['project_resource_summary'].values) # fit has to happen only on train data

# we use the fitted CountVectorizer to convert the text to vector
X_train_prjResSum_ohe = vectorizer.transform(X_train['project_resource_summary'].values)
#X_cv_prjResSum_ohe = vectorizer.transform(X_cv['project_resource_summary'].values)
X_test_prjResSum_ohe = vectorizer.transform(X_test['project_resource_summary'].values)

print("project_resource_summary After vectorizations")
print(X_train_prjResSum_ohe.shape, y_train.shape)
#print(X_cv_prjResSum_ohe.shape, y_cv.shape)
print(X_test_prjResSum_ohe.shape, y_test.shape)
#print(vectorizer.get_feature_names())
print("="*100)
ff=vectorizer.get_feature_names()
```

```
project_resource_summary After vectorizations
(33500, 10613) (33500,)
(16500, 10613) (16500,)
=====
```

## 2.2 Make Data Model Ready: encoding numerical, categorical features

In [54]:

```
# please write all the code with proper documentation, and proper titles for each subsection
# go through documentations and blogs before you start coding
# first figure out what to do, and then think about how to do.
# reading and understanding error messages will be very much helpfull in debugging your code
# make sure you featurize train and test data separatly

# when you plot any graph make sure you use
# a. Title, that describes your plot, this will be very helpful to the reader
# b. Legends if needed
# c. X-axis label
# d. Y-axis label
```

### 2.2.1 Make Data Model Ready: encoding numerical | quantity

In [55]:

```
from sklearn.preprocessing import Normalizer
normalizer = Normalizer()
# normalizer.fit(X_train['price'].values)
# this will rise an error Expected 2D array, got 1D array instead:
# array=[105.22 215.96 96.01 ... 368.98 80.53 709.67].
# Reshape your data either using
# array.reshape(-1, 1) if your data has a single feature
# array.reshape(1, -1) if it contains a single sample.
normalizer.fit(X_train['quantity'].values.reshape(-1,1))

X_train_quantity_norm = normalizer.transform(X_train['quantity'].values.reshape(-1,1))
#X_cv_quantity_norm = normalizer.transform(X_cv['quantity'].values.reshape(-1,1))
X_test_quantity_norm = normalizer.transform(X_test['quantity'].values.reshape(-1,1))

print("quantity After vectorizations")
print(X_train_quantity_norm.shape, y_train.shape)
#print(X_cv_quantity_norm.shape, y_cv.shape)
print(X_test_quantity_norm.shape, y_test.shape)
print("=="*100)
```

```
quantity After vectorizations
(33500, 1) (33500,)
(16500, 1) (16500,)
=====
```

### 2.2.2 Make Data Model Ready: encoding numerical| teacher\_number\_of\_previously\_posted\_projects

In [56]:

```
from sklearn.preprocessing import Normalizer
normalizer = Normalizer()
# normalizer.fit(X_train['price'].values)
# this will rise an error Expected 2D array, got 1D array instead:
# array=[105.22 215.96 96.01 ... 368.98 80.53 709.67].
# Reshape your data either using
# array.reshape(-1, 1) if your data has a single feature
# array.reshape(1, -1) if it contains a single sample.
```

```
normalizer.fit(X_train['teacher_number_of_previously_posted_projects'].values.reshape(-1,1))

X_train_TprevPrj_norm =
normalizer.transform(X_train['teacher_number_of_previously_posted_projects'].values.reshape(-1,1))
#X_cv_TprevPrj_norm =
normalizer.transform(X_cv['teacher_number_of_previously_posted_projects'].values.reshape(-1,1))
X_test_TprevPrj_norm = normalizer.transform(X_test['teacher_number_of_previously_posted_projects']
.values.reshape(-1,1))

print("teacher_number_of_previously_posted_projects After vectorizations")
print(X_train_TprevPrj_norm.shape, y_train.shape)
#print(X_cv_TprevPrj_norm.shape, y_cv.shape)
print(X_test_TprevPrj_norm.shape, y_test.shape)
print("="*100)
```

```
teacher_number_of_previously_posted_projects After vectorizations
(33500, 1) (33500,)
(16500, 1) (16500,)
=====
```

## 2.2.3 Make Data Model Ready: encoding numerical | price

In [57]:

```
from sklearn.preprocessing import Normalizer
normalizer = Normalizer()
# normalizer.fit(X_train['price'].values)
# this will rise an error Expected 2D array, got 1D array instead:
# array=[105.22 215.96 96.01 ... 368.98 80.53 709.67].
# Reshape your data either using
# array.reshape(-1, 1) if your data has a single feature
# array.reshape(1, -1) if it contains a single sample.
normalizer.fit(X_train['price'].values.reshape(-1,1))

X_train_price_norm = normalizer.transform(X_train['price'].values.reshape(-1,1))
#X_cv_price_norm = normalizer.transform(X_cv['price'].values.reshape(-1,1))
X_test_price_norm = normalizer.transform(X_test['price'].values.reshape(-1,1))

print("Price After vectorizations")
print(X_train_price_norm.shape, y_train.shape)
#print(X_cv_price_norm.shape, y_cv.shape)
print(X_test_price_norm.shape, y_test.shape)
print("="*100)
```

```
Price After vectorizations
(33500, 1) (33500,)
(16500, 1) (16500,)
=====
```

In [58]:

```
h=['price','quantity','teacher_number_of_previously_posted_projects']
print(type(h))
```

```
<class 'list'>
```

## 2.2.4 Make Data Model Ready: encoding numerical | sentimental score

### 2.2.4.1 Make Data Model Ready: encoding numerical | sentimental score | neg

In [59]:

```
from sklearn.preprocessing import Normalizer
normalizer = Normalizer()
# normalizer.fit(X_train['price'].values)
# this will rise an error Expected 2D array, got 1D array instead:
```



```
# array=[105.22 215.96 96.01 ... 368.98 80.53 709.67].
# Reshape your data either using
# array.reshape(-1, 1) if your data has a single feature
# array.reshape(1, -1) if it contains a single sample.
normalizer.fit(X_train['neg'].values.reshape(-1,1))

X_train_neg_norm = normalizer.transform(X_train['neg'].values.reshape(-1,1))
#X_cv_neg_norm = normalizer.transform(X_cv['neg'].values.reshape(-1,1))
X_test_neg_norm = normalizer.transform(X_test['neg'].values.reshape(-1,1))

print("neg After vectorizations")
print(X_train_neg_norm.shape, y_train.shape)
#print(X_cv_neg_norm.shape, y_cv.shape)
print(X_test_neg_norm.shape, y_test.shape)
print("=="*100)
```

```
neg After vectorizations
(33500, 1) (33500,)
(16500, 1) (16500,)
=====
```

#### 2.2.4.2 Make Data Model Ready: encoding numerical | sentimental score | pos

In [60]:

```
from sklearn.preprocessing import Normalizer
normalizer = Normalizer()
# normalizer.fit(X_train['price'].values)
# this will rise an error Expected 2D array, got 1D array instead:
# array=[105.22 215.96 96.01 ... 368.98 80.53 709.67].
# Reshape your data either using
# array.reshape(-1, 1) if your data has a single feature
# array.reshape(1, -1) if it contains a single sample.
normalizer.fit(X_train['pos'].values.reshape(-1,1))

X_train_pos_norm = normalizer.transform(X_train['pos'].values.reshape(-1,1))
#X_cv_pos_norm = normalizer.transform(X_cv['pos'].values.reshape(-1,1))
X_test_pos_norm = normalizer.transform(X_test['pos'].values.reshape(-1,1))

print("pos After vectorizations")
print(X_train_pos_norm.shape, y_train.shape)
#print(X_cv_pos_norm.shape, y_cv.shape)
print(X_test_pos_norm.shape, y_test.shape)
print("=="*100)
```

```
pos After vectorizations
(33500, 1) (33500,)
(16500, 1) (16500,)
=====
```

#### 2.2.4.3 Make Data Model Ready: encoding numerical | sentimental score | neu

In [61]:

```
from sklearn.preprocessing import Normalizer
normalizer = Normalizer()
# normalizer.fit(X_train['price'].values)
# this will rise an error Expected 2D array, got 1D array instead:
# array=[105.22 215.96 96.01 ... 368.98 80.53 709.67].
# Reshape your data either using
# array.reshape(-1, 1) if your data has a single feature
# array.reshape(1, -1) if it contains a single sample.
normalizer.fit(X_train['neu'].values.reshape(-1,1))

X_train_neu_norm = normalizer.transform(X_train['neu'].values.reshape(-1,1))
#X_cv_neu_norm = normalizer.transform(X_cv['neu'].values.reshape(-1,1))
X_test_neu_norm = normalizer.transform(X_test['neu'].values.reshape(-1,1))

print("neu After vectorizations")
```

```
print(X_train_neu_norm.shape, y_train.shape)
#print(X_cv_neu_norm.shape, y_cv.shape)
print(X_test_neu_norm.shape, y_test.shape)
print("="*100)
```

```
neu After vectorizations
(33500, 1) (33500,)
(16500, 1) (16500,)
=====
```

#### 2.2.4.4 Make Data Model Ready: encoding numerical | sentimental score | compound

In [62]:

```
from sklearn.preprocessing import Normalizer
normalizer = Normalizer()
# normalizer.fit(X_train['price'].values)
# this will rise an error Expected 2D array, got 1D array instead:
# array=[105.22 215.96 96.01 ... 368.98 80.53 709.67].
# Reshape your data either using
# array.reshape(-1, 1) if your data has a single feature
# array.reshape(1, -1) if it contains a single sample.
normalizer.fit(X_train['compound'].values.reshape(-1,1))

X_train_compound_norm = normalizer.transform(X_train['compound'].values.reshape(-1,1))
#X_cv_compound_norm = normalizer.transform(X_cv['compound'].values.reshape(-1,1))
X_test_compound_norm = normalizer.transform(X_test['compound'].values.reshape(-1,1))

print("compound After vectorizations")
print(X_train_compound_norm.shape, y_train.shape)
#print(X_cv_compound_norm.shape, y_cv.shape)
print(X_test_compound_norm.shape, y_test.shape)
print("="*100)
```

```
compound After vectorizations
(33500, 1) (33500,)
(16500, 1) (16500,)
=====
```

#### 2.2.5 Make Data Model Ready: encoding numerical | number of words in the title

In [63]:

```
from sklearn.preprocessing import Normalizer
normalizer = Normalizer()
# normalizer.fit(X_train['price'].values)
# this will rise an error Expected 2D array, got 1D array instead:
# array=[105.22 215.96 96.01 ... 368.98 80.53 709.67].
# Reshape your data either using
# array.reshape(-1, 1) if your data has a single feature
# array.reshape(1, -1) if it contains a single sample.
normalizer.fit(X_train['title_wc'].values.reshape(-1,1))

X_train_title_wc_norm = normalizer.transform(X_train['title_wc'].values.reshape(-1,1))
#X_cv_title_wc_norm = normalizer.transform(X_cv['title_wc'].values.reshape(-1,1))
X_test_title_wc_norm = normalizer.transform(X_test['title_wc'].values.reshape(-1,1))

print("title_wc After vectorizations")
print(X_train_title_wc_norm.shape, y_train.shape)
#print(X_cv_title_wc_norm.shape, y_cv.shape)
print(X_test_title_wc_norm.shape, y_test.shape)
print("="*100)
```

```
title_wc After vectorizations
(33500, 1) (33500,)
(16500, 1) (16500,)
=====
```

## 2.2.6 Make Data Model Ready: encoding numerical | number of words in the essay

In [64]:

```
from sklearn.preprocessing import Normalizer
normalizer = Normalizer()
# normalizer.fit(X_train['price'].values)
# this will rise an error Expected 2D array, got 1D array instead:
# array=[105.22 215.96 96.01 ... 368.98 80.53 709.67].
# Reshape your data either using
# array.reshape(-1, 1) if your data has a single feature
# array.reshape(1, -1) if it contains a single sample.
normalizer.fit(X_train['essay_wc'].values.reshape(-1,1))

X_train_essay_wc_norm = normalizer.transform(X_train['essay_wc'].values.reshape(-1,1))
#X_cv_essay_wc_norm = normalizer.transform(X_cv['essay_wc'].values.reshape(-1,1))
X_test_essay_wc_norm = normalizer.transform(X_test['essay_wc'].values.reshape(-1,1))

print("essay_wc After vectorizations")
print(X_train_essay_wc_norm.shape, y_train.shape)
#print(X_cv_essay_wc_norm.shape, y_cv.shape)
print(X_test_essay_wc_norm.shape, y_test.shape)
print("=="*100)
```

```
essay_wc After vectorizations
(33500, 1) (33500,)
(16500, 1) (16500,)
```

=====

## 2.3 Make Data Model Ready: encoding eassay, and project\_title

In [65]:

```
# please write all the code with proper documentation, and proper titles for each subsection
# go through documentations and blogs before you start coding
# first figure out what to do, and then think about how to do.
# reading and understanding error messages will be very much helpfull in debugging your code
# make sure you featurize train and test data separatly

# when you plot any graph make sure you use
# a. Title, that describes your plot, this will be very helpful to the reader
# b. Legends if needed
# c. X-axis label
# d. Y-axis label
```

### 2.3.1 Make Data Model Ready: project\_essay | BOW

In [111]:

```
from sklearn.feature_extraction.text import CountVectorizer
# categorical, numerical features + project_title(BOW) + preprocessed_eassay
# (BOW with bi-grams with min_df=10 and max_features=5000)
vectorizer = CountVectorizer(min_df=10,ngram_range=(1,2), max_features=5000)
vectorizer.fit(X_train['essay'].values) # fit has to happen only on train data

# we use the fitted CountVectorizer to convert the text to vector
X_train_essay_bow = vectorizer.transform(X_train['essay'].values)
#X_cv_essay_bow = vectorizer.transform(X_cv['essay'].values)
X_test_essay_bow = vectorizer.transform(X_test['essay'].values)

print("Essay After vectorizations")
print(X_train_essay_bow.shape, y_train.shape)
```

```

#print(X_cv_essay_bow.shape, y_cv.shape)
print(X_test_essay_bow.shape, y_test.shape)
print("="*100)
g=vectorizer.get_feature_names()

```

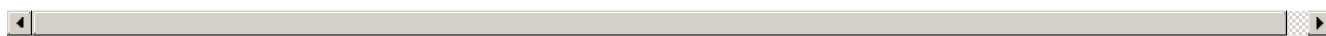
Essay After vectorizations

```

(33500, 5000) (33500,)
(16500, 5000) (16500,)

```

=====



## 2.3.2 Make Data Model Ready: project\_title | BOW

In [67]:

```

vectorizer = CountVectorizer()
# categorical, numerical features + project_title(BOW) + preprocessed_essay
# (BOW with bi-grams with min_df=10 and max_features=5000)
vectorizer = CountVectorizer(min_df=10,ngram_range=(1,2), max_features=5000)
vectorizer.fit(X_train['project_title'].values) # fit has to happen only on train data

# we use the fitted CountVectorizer to convert the text to vector
X_train_title_bow = vectorizer.transform(X_train['project_title'].values)
#X_cv_title_bow = vectorizer.transform(X_cv['project_title'].values)
X_test_title_bow = vectorizer.transform(X_test['project_title'].values)

print("project_title After vectorizations")
print(X_train_title_bow.shape, y_train.shape)
#print(X_cv_title_bow.shape, y_cv.shape)
print(X_test_title_bow.shape, y_test.shape)
#print(vectorizer.get_feature_names())
print("="*100)
k=vectorizer.get_feature_names()

```

project\_title After vectorizations

```

(33500, 3404) (33500,)
(16500, 3404) (16500,)

```

=====



## 2.3.3 Make Data Model Ready: project\_essay | TFIDF

In [68]:

```

from sklearn.feature_extraction.text import TfidfVectorizer
# categorical, numerical features + project_title(BOW) + preprocessed_essay
# (TFIDF with bi-grams with min_df=10 and max_features=5000)
Tfidf_vectorizer = TfidfVectorizer(min_df=10,ngram_range=(1,2), max_features=5000)

Tfidf_vectorizer.fit(X_train['essay'].values)

X_train_text_tfidf = Tfidf_vectorizer.transform(X_train['essay'].values)
#X_cv_text_tfidf = Tfidf_vectorizer.transform(X_cv['essay'].values)
X_test_text_tfidf = Tfidf_vectorizer.transform(X_test['essay'].values)

##print("Shape of matrix after one hot encodig ",text_tfidf.shape)

print("Essay After vectorizations")
print(X_train_text_tfidf.shape, y_train.shape)
#print(X_cv_text_tfidf.shape, y_cv.shape)
print(X_test_text_tfidf.shape, y_test.shape)
#print(Tfidf_vectorizer.get_feature_names())
print("="*100)
ii=Tfidf_vectorizer.get_feature_names()

```

Essay After vectorizations

```

(33500, 5000) (33500,)
(16500, 5000) (16500,)

```

=====

## 2.3.4 Make Data Model Ready: project\_title | TFIDF

In [69]:

```
from sklearn.feature_extraction.text import TfidfVectorizer
# categorical, numerical features + project_title(BOW) + preprocessed_eassay
# (TFIDF with bi-grams with min_df=10 and max_features=5000)
Tfidf_vectorizer = TfidfVectorizer(min_df=10, ngram_range=(1,2), max_features=5000)

Tfidf_vectorizer.fit(X_train['project_title'].values)

X_train_title_tfidf = Tfidf_vectorizer.transform(X_train['project_title'].values)
#X_cv_title_tfidf = Tfidf_vectorizer.transform(X_cv['project_title'].values)
X_test_title_tfidf = Tfidf_vectorizer.transform(X_test['project_title'].values)

##print("Shape of matrix after one hot encodig ",text_tfidf.shape)

print("project_title After vectorizations")
print(X_train_title_tfidf.shape, y_train.shape)
#print(X_cv_title_tfidf.shape, y_cv.shape)
print(X_test_title_tfidf.shape, y_test.shape)
#print(Tfidf_vectorizer.get_feature_names())
print("="*100)
j=Tfidf_vectorizer.get_feature_names()
```

```
project_title After vectorizations
(33500, 3404) (33500,)
(16500, 3404) (16500,)
```

## 2.3.5 Make Data Model Ready: project\_essay | AVG W2V

In [70]:

```
# average Word2Vec for Train Essay
# compute average word2vec for each review.
X_train_essay_avg_w2v = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_train['essay'].values): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if word in glove_words:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    X_train_essay_avg_w2v.append(vector)

print(len(X_train_essay_avg_w2v))
print(len(X_train_essay_avg_w2v[0]))

# stronging variables into pickle files python: http://www.jessicayung.com/how-to-use-pickle-to-sa
ve-and-load-variables-in-python/

import pickle
with open('X_train_essay_avg_w2v', 'wb') as f:
    pickle.dump(X_train_essay_avg_w2v, f)
```

```
100%|██| 33500/33500
[00:12<00:00, 2700.93it/s]
```

```
33500
300
```

In [71]:

```
## stronging variables into pickle files python: http://www.jessicayung.com/how-to-use-pickle-to-s
```

```
ave-and-load-variables-in-python/
## make sure you have the glove_vectors file
#with open#('X_train_essay_avg_w2v', 'rb') as f:
#    X_train_essay_avg_w2v = pickle.load(f)
#
#print(len(X_train_essay_avg_w2v))
#print(len(X_train_essay_avg_w2v[0]))
```

In [72]:

```
# average Word2Vec for Test Essay
# compute average word2vec for each review.
X_test_essay_avg_w2v = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_test['essay'].values): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if word in glove_words:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    X_test_essay_avg_w2v.append(vector)

print(len(X_test_essay_avg_w2v))
print(len(X_test_essay_avg_w2v[0]))
```

```
100%|██| 16500/16500
[00:06<00:00, 2733.57it/s]
```

```
16500
300
```

## 2.3.6 Make Data Model Ready: project\_title | AVG W2V

In [73]:

```
# average Word2Vec for Train Title
# compute average word2vec for each review.
X_train_title_avg_w2v = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_train['project_title'].values): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if word in glove_words:
            vector += model[word]
            cnt_words += 1
    if cnt_words != 0:
        vector /= cnt_words
    X_train_title_avg_w2v.append(vector)

print(len(X_train_title_avg_w2v))
print(len(X_train_title_avg_w2v[0]))
```

```
100%|██| 33500/33500
[00:00<00:00, 126261.46it/s]
```

```
33500
300
```

In [74]:

```
# average Word2Vec for Test Essay
# compute average word2vec for each review.
X_test_title_avg_w2v = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_test['project_title'].values): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    cnt_words = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if word in glove_words:
            vector += model[word]
```

```
100%|██| 16500/16500  
[00:00<00:00, 125334.16it/s]
```

```
100%|██| 33500/33500 [02:  
15<00:00, 247.93it/s]
```

```
# TFIDF weighted Word2Vec for test essay
```

```
# compute average word2vec for each review.
te_tfidf_w2v_essay_vectors = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_test['essay'].values): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in te_tfidf_model_essay):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf
            value((sentence.count(word)/len(sentence.split())))
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tf
            idf value for each word
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    te_tfidf_w2v_essay_vectors.append(vector)

print(len(te_tfidf_w2v_essay_vectors))
print(len(te_tfidf_w2v_essay_vectors[0]))
```

```
100%|███| 16500/16500 [01:  
05<00:00, 250.23it/s]
```

16500  
300

### 2.3.8 Make Data Model Ready: project\_title | TFIDF W2V

In [79]:

```
# S = ["abc def pqr", "def def def abc", "pqr pqr def"]
Tr_tfidf_model_title = TfidfVectorizer()
Tr_tfidf_model_title.fit(X_train['project_title'].values)
# we are converting a dictionary with word as a key, and the idf as a value
dictionary = dict(zip(Tr_tfidf_model_title.get_feature_names(), list(Tr_tfidf_model_title.idf_)))
m=Tr_tfidf_model_title.get_feature_names()
Tr_tfidf_model_title = set(Tr_tfidf_model_title.get_feature_names())
```

In [80]:

```
# TFIDF weighted Word2Vec for train title
# compute average word2vec for each review.
tr_tfidf_w2v_title_vectors = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_train['project_title'].values): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in Tr_tfidf_model_title):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf
            value((sentence.count(word)/len(sentence.split()))))
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tf
            idf value for each word
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tr_tfidf_w2v_title_vectors.append(vector)

print(len(tr_tfidf_w2v_title_vectors))
print(len(tr_tfidf_w2v_title_vectors[0]))
```

```
100%|██| 33500/33500  
[00:00<00:00, 93554.85it/s]
```

33500  
300



```
# S = ["abc def pqr", "def def def abc", "pqr pqr def"]
te_tfidf_model_title = TfidfVectorizer()
te_tfidf_model_title.fit(X_test['project_title'].values)
# we are converting a dictionary with word as a key, and the idf as a value
dictionary = dict(zip(te_tfidf_model_title.get_feature_names(), list(te_tfidf_model_title.idf_)))
mm=te_tfidf_model_title.get_feature_names()
te_tfidf_model_title = set(te_tfidf_model_title.get_feature_names())
```

```
# TFIDF weighted Word2Vec for test title
# compute average word2vec for each review.
te_tfidf_w2v_title_vectors = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(X_test['project_title'].values): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in te_tfidf_model_title):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf
            value((sentence.count(word)/len(sentence.split())))
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tf
            idf value for each word
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    te_tfidf_w2v_title_vectors.append(vector)

print(len(te_tfidf_w2v_title_vectors))
print(len(te_tfidf_w2v_title_vectors[0]))
```

16500  
300

Apply Decision Tree on different kind of featurization as mentioned in the instructions  
For Every model that you work on make sure you do the step 2 and step 3 of instructions

```
# please write all the code with proper documentation, and proper titles for each subsection
# go through documentations and blogs before you start coding
# first figure out what to do, and then think about how to do.
# reading and understanding error messages will be very much helpfull in debugging your code
# when you plot any graph make sure you use
    # a. Title, that describes your plot, this will be very helpful to the reader
    # b. Legends if needed
    # c. X-axis label
    # d. Y-axis label
```

1. Apply Decision Tree Classifier(DecisionTreeClassifier) on these feature sets

- ```
# Please write all the code with proper documentation
```

In [85]:

```
# merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
from scipy.sparse import hstack
X_tr_bow = hstack((X_train_essay_bow, X_train_title_bow, X_train_state_oh, X_train_clean_oh,
X_train_cleanSub_oh, X_train_grade_oh, X_train_teacher_oh, X_train_prjResSum_oh,
X_train_quantity_norm, X_train_TprevPrj_norm, X_train_price_norm)).tocsr()
X_te_bow = hstack((X_test_essay_bow, X_test_title_bow, X_test_state_oh, X_test_clean_oh, X_test_
cleanSub_oh, X_test_grade_oh, X_test_teacher_oh, X_test_prjResSum_oh, X_test_quantity_norm,
X_test_TprevPrj_norm, X_test_price_norm)).tocsr()

print("Final Data matrix | BOW")
print(X_tr_bow.shape, y_train.shape)
print(X_te_bow.shape, y_test.shape)
print("=="*100)
```

```
Final Data matrix | BOW
(33500, 19119) (33500,)
(16500, 19119) (16500,)
```

In [86]:

```
#code source:
http://occam.olin.edu/sites/default/files/DataScienceMaterials/machine_learning_lecture_2/Machine%
rning%20Lecture%202.html
from sklearn.model_selection import train_test_split
#from sklearn.grid_search import GridSearchCV
from sklearn.model_selection import GridSearchCV
from sklearn.datasets import *
from sklearn.tree import DecisionTreeClassifier

d_range=[1,5,10,50,100,500,1000]
split_range=[5,10,100,500]

param_grid=dict(max_depth=d_range,min_samples_split=split_range)
print(param_grid)

#Using GridSearchCV
# https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.GridSearchCV.html
modelBow = GridSearchCV(DecisionTreeClassifier(class_weight="balanced"), param_grid, scoring =
'f1', cv=5)
modelBow.fit(X_tr_bow, y_train)

print(modelBow.best_estimator_)
print(modelBow.score(X_te_bow, y_test))
```

```
{'max_depth': [1, 5, 10, 50, 100, 500, 1000], 'min_samples_split': [5, 10, 100, 500]}
DecisionTreeClassifier(class_weight='balanced', criterion='gini', max_depth=1,
                        max_features=None, max_leaf_nodes=None,
                        min_impurity_decrease=0.0, min_impurity_split=None,
                        min_samples_leaf=1, min_samples_split=5,
                        min_weight_fraction_leaf=0.0, presort=False, random_state=None,
                        splitter='best')
0.8976886471787899
```

**"1. When you consider AUC as a metric, just plot AUC vs max\_depth. You have to plot two curves, one curve from training data and another from cross-validation data in same plot. So that you will get a clear idea of when a model is overfitting or underfitting. You have to choose the hyperparameter before it overfits or underfits. Same can be followed for AUC vs min\_sample\_split.(Please sort the AUC values and plot if graph you are getting is varying too much) OR you can plot heatmaps or 3D-plots, you can find code from 3d\_scatter\_plot.ipynb (Please upload the screen shots of 3D plots).**

In [88]:

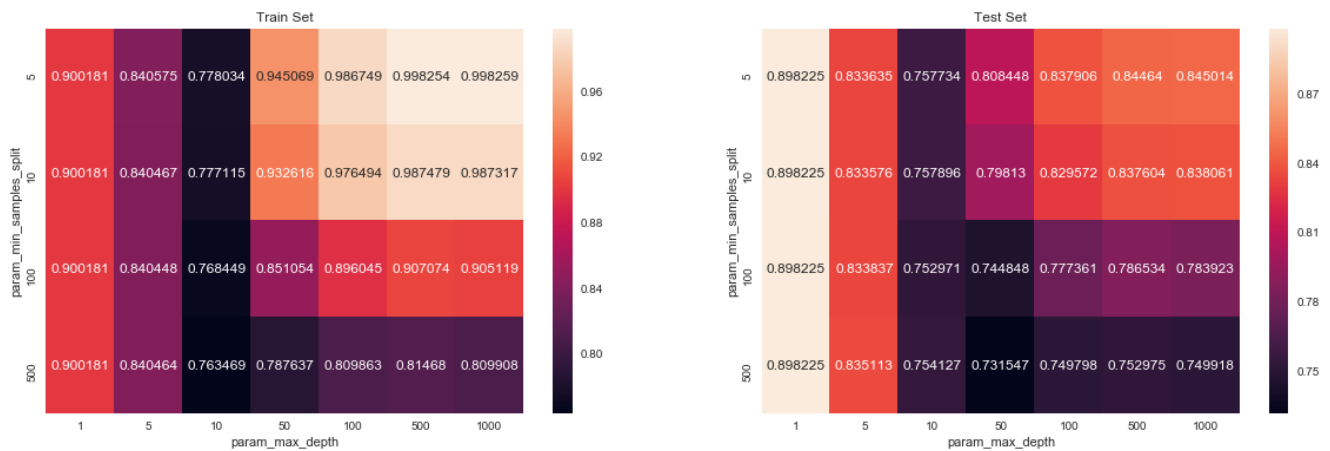
```
# https://seaborn.pydata.org/generated/seaborn.heatmap.html
import seaborn as sns; sns.set()
max_scores1=pd.DataFrame(modelBow.cv_results_).groupby(['param_min_samples_split','param_max_depth'
]).max().unstack()[['mean_test_score','mean_train_score']]
```

```
fig,ax=plt.subplots(1,2,figsize=(20,6))

sns.heatmap(max_scores1.mean_train_score,annot=True,fmt='4g',ax=ax[0])
sns.heatmap(max_scores1.mean_test_score,annot=True,fmt='4g',ax=ax[1])

ax[0].set_title('Train Set')
ax[1].set_title('Test Set')

plt.show()
```



## Conclusion

For all the various values of min\_samples\_split=5 and max\_depth=1 is giving the best score for test data.

In [89]:

```
#code source:
http://occam.olin.edu/sites/default/files/DataScienceMaterials/machine_learning_lecture_2/Machine%20Learning%20Lecture%20202.html
from sklearn.model_selection import train_test_split
#from sklearn.grid_search import GridSearchCV
from sklearn.model_selection import GridSearchCV
from sklearn.datasets import *
from sklearn.tree import DecisionTreeClassifier

d_range=[1]
split_range=[5]

param_grid=dict(max_depth=d_range,min_samples_split=split_range)
print(param_grid)

#Using GridSearchCV
# https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.GridSearchCV.html
modelBowB = GridSearchCV(DecisionTreeClassifier(class_weight="balanced"), param_grid, scoring = 'f1', cv=5)
modelBowB.fit(X_tr_bow, y_train)

print(modelBowB.best_estimator_)
print(modelBowB.score(X_te_bow, y_test))

{'max_depth': [1], 'min_samples_split': [5]}
DecisionTreeClassifier(class_weight='balanced', criterion='gini', max_depth=1,
                        max_features=None, max_leaf_nodes=None,
                        min_impurity_decrease=0.0, min_impurity_split=None,
                        min_samples_leaf=1, min_samples_split=5,
                        min_weight_fraction_leaf=0.0, presort=False, random_state=None,
                        splitter='best')
0.8976886471787899
```

In [90]:

```
best_tuned_parameters = [{'max_depth': [1], 'min_samples_split': [5]}]
```

In [91]:

```
# https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve
from sklearn.metrics import roc_curve, auc

#model = GridSearchCV(LogisticRegression(), best_tuned_parameters)
modelBowB = GridSearchCV(DecisionTreeClassifier(), best_tuned_parameters)
modelBowB.fit(X_tr_bow, y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
# not the predicted outputs
#print(type(model.predict_proba(X_tr_bow)))
#print(model.predict_proba(X_tr_bow))
#print(model.predict_proba(X_tr_bow)[:,-1])

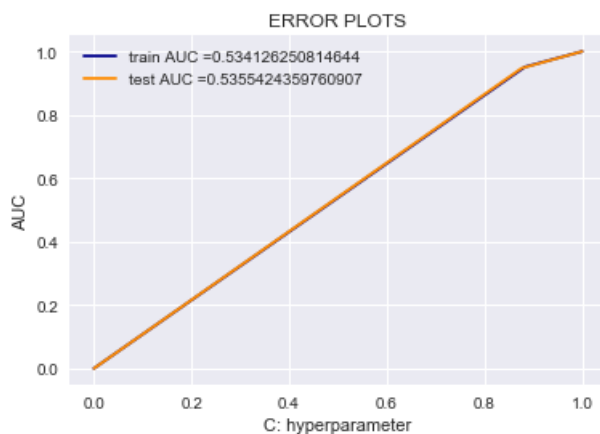
y_train_bow_pred = modelBowB.predict_proba(X_tr_bow)[:,-1]
y_test_bow_pred = modelBowB.predict_proba(X_te_bow)[:,-1]

print(modelBowB.best_estimator_)
print(modelBowB.score(X_te_bow, y_test))

train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_bow_pred)
test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_bow_pred)

plt.plot(train_fpr, train_tpr, label="train AUC =" + str(auc(train_fpr, train_tpr)), color='darkblue')
plt.plot(test_fpr, test_tpr, label="test AUC =" + str(auc(test_fpr, test_tpr)), color='darkorange')
plt.legend()
plt.xlabel("C: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid(True)
plt.show()
```

```
DecisionTreeClassifier(class_weight=None, criterion='gini', max_depth=1,
                        max_features=None, max_leaf_nodes=None,
                        min_impurity_decrease=0.0, min_impurity_split=None,
                        min_samples_leaf=1, min_samples_split=5,
                        min_weight_fraction_leaf=0.0, presort=False, random_state=None,
                        splitter='best')
0.8456969696969697
```



In [92]:

```
# we are writing our own function for predict, with defined threshold
# we will pick a threshold that will give the least fpr
def predict(proba, threshold, fpr, tpr):

    t = threshold[np.argmax(tpr*(1-fpr))]

    # (tpr*(1-fpr)) will be maximum if your fpr is very low and tpr is very high

    print("the maximum value of tpr*(1-fpr)", max(tpr*(1-fpr)), "for threshold", np.round(t,3))
    predictions = []
    for i in proba:
        if i>=t:
            predictions.append(1)
    ,
```

```

else:
    predictions.append(0)
return predictions

```

In [93]:

```

# https://scikit-learn.org/stable/modules/generated/sklearn.metrics.confusion_matrix.html
print("="*100)
from sklearn.metrics import confusion_matrix
print("Train confusion matrix")
print(confusion_matrix(y_train, predict(y_train_bow_pred, tr_thresholds, train_fpr, train_tpr)))
print("Test confusion matrix")
print(confusion_matrix(y_test, predict(y_test_bow_pred, te_thresholds, test_fpr, test_tpr)))

```

```

=====

Train confusion matrix
the maximum value of tpr*(1-fpr) 0.11135206899920401 for threshold 0.855
[[ 605  4563]
 [ 1383 26949]]
Test confusion matrix
the maximum value of tpr*(1-fpr) 0.1175342031122122 for threshold 0.855
[[ 316  2230]
 [ 740 13214]]

```

In [94]:

```

import seaborn as snTr
import seaborn as snTe
import pandas as pdH
import matplotlib.pyplot as pltTr
import matplotlib.pyplot as pltTe

# https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
arrayTr=confusion_matrix(y_train, predict(y_train_bow_pred, tr_thresholds, train_fpr, train_tpr))
df_cmTr = pdH.DataFrame(arrayTr,range(2),range(2))
#print(arrayTr)
# https://stackoverflow.com/questions/32723798/how-do-i-add-a-title-to-seaborn-heatmap
axTr = pltTr.axes()

snTr.set(font_scale=1.4)#for label size

# https://seaborn.pydata.org/generated/seaborn.heatmap.html

snTr.heatmap(df_cmTr, annot=True,annot_kws={"size": 12},fmt="d",ax=axTr)# font size, format in
digit

labels=['Not Approved','Approved']
axTr.set_xticklabels(labels)
axTr.set_yticklabels(labels)
#Suggestion 4.Label confusion matrix heatmap with actual and predicted labels.
pltTr.title("Train confusion matrix")
pltTr.xlabel("Predicted")
pltTr.ylabel("Actual")
pltTr.show()

# https://stackoverflow.com/questions/50947776/plot-two-seaborn-heatmap-graphs-side-by-side
#fig, ax =plt.subplots(1,1)

# https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
arrayTe=confusion_matrix(y_test, predict(y_test_bow_pred, te_thresholds, test_fpr, test_tpr))
df_cmTe = pdH.DataFrame(arrayTe,range(2),range(2))

axTe = pltTe.axes()

snTe.set(font_scale=1.4)#for label size

# https://seaborn.pydata.org/generated/seaborn.heatmap.html
snTe.heatmap(df_cmTe, annot=True,annot_kws={"size": 12},fmt="d",ax=axTe)# font size, format in
digit

#Suggestion 4.Label confusion matrix heatmap with actual and predicted labels.
axTe.set_xticklabels(labels)
axTe.set_yticklabels(labels)
pltTe.title("Test confusion matrix")

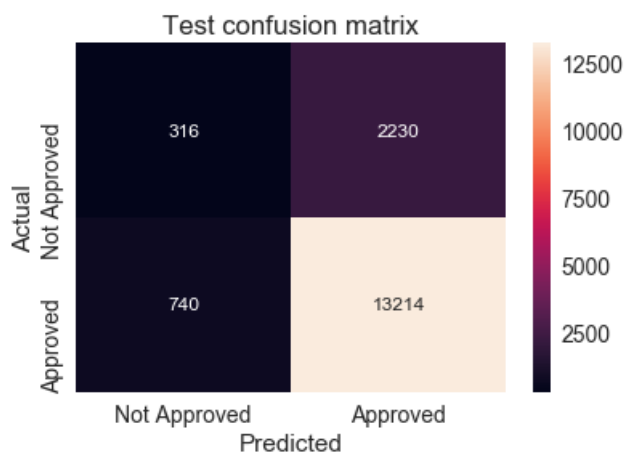
```

```
pltTe.xlabel("Predicted")
pltTe.ylabel("Actual")
pltTe.show()
```

the maximum value of  $tpr \cdot (1-fpr)$  0.11135206899920401 for threshold 0.855



the maximum value of  $tpr \cdot (1-fpr)$  0.1175342031122122 for threshold 0.855



1. Once after you plot the confusion matrix with the test data, get all the 'false positive data points'

- Plot the WordCloud [WordCloud](#)
- Plot the box plot with the 'price' of these 'false positive data points'
- Plot the pdf with the 'teacher\_number\_of\_previously\_posted\_projects' of these 'false positive data points'

In [95]:

```
predict_bow = predict(y_test_bow_pred, te_thresholds, test_fpr, test_tpr) # <=<==

#converting predict_bow list into numpy array.
# https://likegeeks.com/numpy-array-tutorial/
import numpy as np
y_predict = np.array(predict_bow)
print(y_predict)
print(type(y_predict))

# traversing y_test and y_predict, and created new list y_FP, which is 1 for FALSE_POSITIVE, and 0
for rest
y_FP=[]
count_FP=0
# https://www.geeksforgeeks.org/numpy-iterating-over-array/
for ac,pre in np.nditer([y_test,y_predict]):
    if(ac==0 and pre==1):
        y_FP.append(1)
        count_FP+=1
    else:
        y_FP.append(0)

print("False Positive count:",count_FP)
```

```

print(type(y_test))
print(y_test)
print(type(y_predict))

print(y_predict)
print(type(y_FP))
#print(y_FP)
#converting predict_bow list into numpy array.
y_FP=np.array(y_FP)
print(type(y_FP))
print(y_FP)

```

the maximum value of tpr\*(1-fpr) 0.1175342031122122 for threshold 0.855

```

[1 1 1 ... 1 1 1]
<class 'numpy.ndarray'>
False Positive count: 2230
<class 'numpy.ndarray'>
[1 1 1 ... 0 0 1]
<class 'numpy.ndarray'>
[1 1 1 ... 1 1 1]
<class 'list'>
<class 'numpy.ndarray'>
[0 0 0 ... 1 1 0]

```

In [96]:

```

#X_test.columns
#X_test#_working.columns

```

In [97]:

```

X_test_working = X_test.copy(deep=True)  #<=<==

#print("X_test_working length:",len(X_test_working))
print(type(X_test))
print(X_test.shape)
print(type(X_test_working))
print(X_test_working.head(2))

# Adding 3 numpy array into dataframe
print(type(y_test))
print(type(y_predict))
print(type(y_FP))
print(y_test)
print(y_predict)
print(y_FP)

print(type(X_test_working))
X_test_working['Y_Actual'] = y_test
print(X_test_working.head(2))

```

```

<class 'pandas.core.frame.DataFrame'>
(16500, 30)
<class 'pandas.core.frame.DataFrame'>
  Unnamed: 0      id      teacher_id school_state \
33081      141937  p034157  0d7b3cd172c5b19f83a0ed303f46b729      AR
26184      33737  p225681  d44f8cb33de45fce2126bb699d302808      TN

  project_submitted_datetime \
33081      2016-09-26 16:20:26
26184      2016-12-26 14:31:27

  project_title \
33081  Scratching the Surface in Collaborative Learning!
26184      Ready 2 Read!

  project_essay_1 \
33081  With the upcoming school year, I am fortunate ...
26184  I am truly blessed to have the opportunity to ...

  project_essay_2 project_essay_3 \
33081  As my students become more technologically inc...      NaN
26184  My classroom consists of students on all learn...      NaN

```

```

project_essay_4      ...      neu  compound  price quantity \
33081      NaN      ...      0.826    0.9944  438.89      2
26184      NaN      ...      0.812    0.9865  232.49      6

essay_wc essay_len title_wc title_len  prj_res_sum_wc  prj_res_sum_len
33081      213      1528      4      40      14      106
26184      107      763      3      12      9      61

[2 rows x 30 columns]
<class 'numpy.ndarray'>
<class 'numpy.ndarray'>
<class 'numpy.ndarray'>
[1 1 1 ... 0 0 1]
[1 1 1 ... 1 1 1]
[0 0 0 ... 1 1 0]
<class 'pandas.core.frame.DataFrame'>
  Unnamed: 0      id      teacher_id school_state \
33081      141937  p034157  0d7b3cd172c5b19f83a0ed303f46b729      AR
26184      33737  p225681  d44f8cb33de45fced2126bb699d302808      TN

project_submitted_datetime \
33081      2016-09-26 16:20:26
26184      2016-12-26 14:31:27

project_title \
33081  Scratching the Surface in Collaborative Learning!
26184      Ready 2 Read!

project_essay_1 \
33081  With the upcoming school year, I am fortunate ...
26184  I am truly blessed to have the opportunity to ...

project_essay_2 project_essay_3 \
33081  As my students become more technologically inc...      NaN
26184  My classroom consists of students on all learn...      NaN

project_essay_4      ...      compound  price  quantity essay_wc essay_len \
33081      NaN      ...      0.9944  438.89      2      213      1528
26184      NaN      ...      0.9865  232.49      6      107      763

title_wc title_len prj_res_sum_wc  prj_res_sum_len  Y_Actual
33081      4      40      14      106      1
26184      3      12      9      61      1

[2 rows x 31 columns]

```

In [98]:

```

# how to add numpy array into dataframe | https://stackoverflow.com/questions/26666919/add-column-
in-dataframe-from-list
print(X_test_working.columns)
X_test_working['Y_Predict'] = y_predict
X_test_working['Y_FP'] = y_FP
print(X_test_working.columns)

```

```

Index(['Unnamed: 0', 'id', 'teacher_id', 'school_state',
      'project_submitted_datetime', 'project_title', 'project_essay_1',
      'project_essay_2', 'project_essay_3', 'project_essay_4',
      'project_resource_summary',
      'teacher_number_of_previously_posted_projects', 'project_is_approved',
      'clean_categories', 'clean_subcategories', 'clean_grade',
      'clean_tea_pfx', 'essay', 'neg', 'pos', 'neu', 'compound', 'price',
      'quantity', 'essay_wc', 'essay_len', 'title_wc', 'title_len',
      'prj_res_sum_wc', 'prj_res_sum_len', 'Y_Actual'],
      dtype='object')
Index(['Unnamed: 0', 'id', 'teacher_id', 'school_state',
      'project_submitted_datetime', 'project_title', 'project_essay_1',
      'project_essay_2', 'project_essay_3', 'project_essay_4',
      'project_resource_summary',
      'teacher_number_of_previously_posted_projects', 'project_is_approved',
      'clean_categories', 'clean_subcategories', 'clean_grade',
      'clean_tea_pfx', 'essay', 'neg', 'pos', 'neu', 'compound', 'price',
      'quantity', 'essay_wc', 'essay_len', 'title_wc', 'title_len',
      'prj_res_sum_wc', 'prj_res_sum_len', 'Y_Actual', 'Y_Predict', 'Y_FP'],
      dtype='object')

```



In [99]:

```
X_test_FP=X_test_working.copy(deep=True)
X_test_FP.teacher_number_of_previously_posted_projects
X_test_FP.price
X_test_FP.columns
```

Out[99]:

```
Index(['Unnamed: 0', 'id', 'teacher_id', 'school_state',
      'project_submitted_datetime', 'project_title', 'project_essay_1',
      'project_essay_2', 'project_essay_3', 'project_essay_4',
      'project_resource_summary',
      'teacher_number_of_previously_posted_projects', 'project_is_approved',
      'clean_categories', 'clean_subcategories', 'clean_grade',
      'clean_tea_pfx', 'essay', 'neg', 'pos', 'neu', 'compound', 'price',
      'quantity', 'essay_wc', 'essay_len', 'title_wc', 'title_len',
      'prj_res_sum_wc', 'prj_res_sum_len', 'Y_Actual', 'Y_Predict', 'Y_FP'],
      dtype='object')
```

In [100]:

```
X_test_FP=X_test_FP.loc[:,["teacher_number_of_previously_posted_projects","essay","price","Y_FP"]]
X_test_FP.head(2)
```

Out[100]:

	teacher_number_of_previously_posted_projects	essay	price	Y_FP
33081	43	With the upcoming school year, I am fortunate ...	438.89	0
26184	1	I am truly blessed to have the opportunity to ...	232.49	0

In [101]:

```
print(X_tr_bow.shape, y_train.shape)
print(X_te_bow.shape, y_test.shape)
```

```
(33500, 19119) (33500,)
(16500, 19119) (16500,)
```

In [102]:

```
X_test_FP = X_test_FP[X_test_FP.Y_FP==True]
```

In [103]:

```
# https://www.geeksforgeeks.org/generating-word-cloud-python/
# Python program to generate WordCloud

# importing all necessary modules
from wordcloud import WordCloud, STOPWORDS
import matplotlib.pyplot as plt
import pandas as pd

# Reads 'Youtube04-Eminem.csv' file
df = pd.read_csv(r"Youtube04-Eminem.csv", encoding="latin-1")

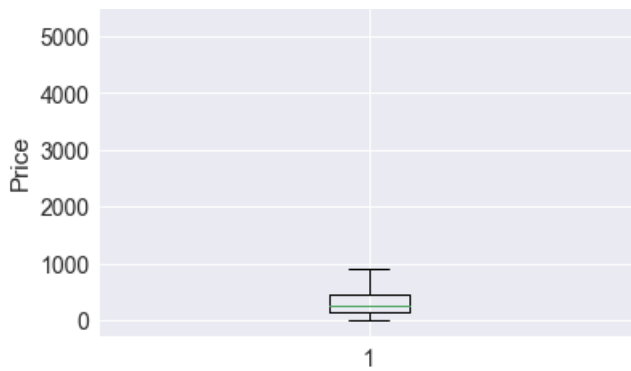
comment_words = ' '
stopwords = set(STOPWORDS)

# iterate through the csv file
for val in X_test_FP["essay"][:1]:

    # typecaste each val to string
    val = str(val)

    # split the value
    tokens = val.split()
```





**Plot the pdf with the teacher\_number\_of\_previously\_posted\_projects of these false positive data points**

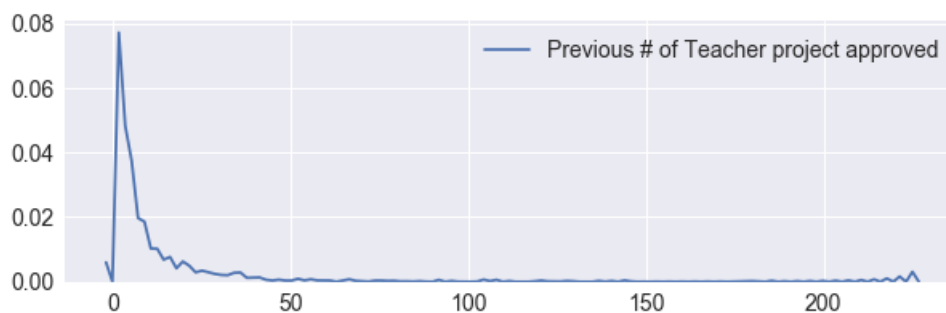
In [105]:

```
plt.figure(figsize=(10,3))

# https://seaborn.pydata.org/generated/seaborn.kdeplot.html | kernel density estimate | sns.kdeplot
# bw : {'scott' | 'silverman' | scalar | pair of scalars }, optional
# Name of reference method to determine kernel size, scalar factor, or scalar for each dimension
# of the bivariate plot. Note that the underlying computational libraries have different
# interpretations
# for this parameter: statsmodels uses it directly, but scipy treats it as a scaling factor
# for the standard deviation of the data.

sns.kdeplot(X_test_FP["teacher_number_of_previously_posted_projects"],label="Previous # of Teacher
project approved", bw=0.6)
plt.legend()
plt.show()

#PDF
```



#### 2.4.1.1 Graphviz visualization of Decision Tree on BOW, SET 1

In [106]:

```
# Please write all the code with proper documentation
```

In [113]:

```
print(len(aa))
print(len(b))
print(len(c))
print(len(d))
print(len(e))
print(len(ff))
print(len(g))
print(len(k))
print(len(h))
#print(len(ii))
#print(len(j))
#print(len(l))
```

```
51
9
30
4
5
10613
5000
3404
3
```

In [114]:

```
feature_list=[]
print(feature_list)
feature_list=aa+b+c+d+e+ff+g+k+h
#print(feature_list)
print(len(feature_list))
```

```
[]
19119
```

In [115]:

```
print(X_tr_bow[:,1].shape)
print(type(X_tr_bow))
print(y_train.shape)
print(type(y_train))
```

```
(33500, 19119)
<class 'scipy.sparse.csr.csr_matrix'>
(33500,)
<class 'numpy.ndarray'>
```

In [116]:

```
#print(feature_list)
```

## Suggestion

"Your graphviz outputs are not available in your work so that is the we asked you to upload them as .png files in suggestion 2 please include it and re-submit your work"

In [119]:

```
# error | NotFittedError: This GridSearchCV instance is not fitted yet. Call 'fit' with
appropriate arguments before using this method.
# https://stackoverflow.com/questions/46192063/not-fitted-error-when-using-sklearns-graphviz

from sklearn import tree
import graphviz as g
clf = tree.DecisionTreeClassifier(max_depth=2, min_samples_split=5, class_weight="balanced")
clf = clf.fit(X_tr_bow, y_train)
import pydot

dot_data = tree.export_graphviz(clf, out_file='X_tr_bow_graphviz.dot', max_depth=2,
                                feature_names=feature_list,
                                class_names=True,
                                filled=True, rounded=True,
                                special_characters=False)

graph = g.Source(dot_data)
graph

# https://scikit-learn.org/stable/modules/generated/sklearn.tree.export_graphviz.html#sklearn.tree.export_graphviz
# https://stackoverflow.com/questions/5316206/convert-dot-to-png-in-python

(graph,) = pydot.graph_from_dot_file('X_tr_bow_graphviz.dot')
graph.write_png('X_tr_bow_graphviz.png')
```

## 2.4.2 Applying Decision Trees on TFIDF, SET 2

### 1. Apply Decision Tree Classifier(DecisionTreeClassifier) on these feature sets

- **Set 2:** categorical, numerical features + project\_title(TFIDF)+ preprocessed\_eassay (TFIDF)

In [120]:

```
# merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
from scipy.sparse import hstack
X_tr_tfidf = hstack((X_train_text_tfidf, X_train_title_tfidf, X_train_state_oh, X_train_clean_oh,
, X_train_cleanSub_oh, X_train_grade_oh, X_train_teacher_oh, X_train_prjResSum_oh,
X_train_quantity_norm, X_train_TprevPrj_norm, X_train_price_norm)).tocsr()
X_te_tfidf = hstack((X_test_text_tfidf, X_test_title_tfidf, X_test_state_oh, X_test_clean_oh, X_
test_cleanSub_oh, X_test_grade_oh, X_test_teacher_oh, X_test_prjResSum_oh,
X_test_quantity_norm, X_test_TprevPrj_norm, X_test_price_norm)).tocsr()

print("Final Data matrix | tfidf")
print(X_tr_tfidf.shape, y_train.shape)
print(X_te_tfidf.shape, y_test.shape)
print("=="*100)
```

```
Final Data matrix | tfidf
(33500, 19119) (33500,)
(16500, 19119) (16500,)
```

In [121]:

```
#code source:
http://occam.olin.edu/sites/default/files/DataScienceMaterials/machine_learning_lecture_2/Machine%
rning%20Lecture%202.html
from sklearn.model_selection import train_test_split
#from sklearn.grid_search import GridSearchCV
from sklearn.model_selection import GridSearchCV
from sklearn.datasets import *
from sklearn.tree import DecisionTreeClassifier

d_range=[1,5,10,50,100,500,1000]
split_range=[5,10,100,500]

param_grid=dict(max_depth=d_range,min_samples_split=split_range)
print(param_grid)

#Using GridSearchCV
# https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.GridSearchCV.html
modeltfidf = GridSearchCV(DecisionTreeClassifier(class_weight="balanced"), param_grid, scoring = '
f1', cv=5)
modeltfidf.fit(X_tr_tfidf, y_train)

print(modeltfidf.best_estimator_)
print(modeltfidf.score(X_te_tfidf, y_test))
```

```
{'max_depth': [1, 5, 10, 50, 100, 500, 1000], 'min_samples_split': [5, 10, 100, 500]}
DecisionTreeClassifier(class_weight='balanced', criterion='gini',
                        max_depth=500, max_features=None, max_leaf_nodes=None,
                        min_impurity_decrease=0.0, min_impurity_split=None,
                        min_samples_leaf=1, min_samples_split=5,
                        min_weight_fraction_leaf=0.0, presort=False, random_state=None,
                        splitter='best')
0.846117441479486
```

**"1. When you consider AUC as a metric, just plot AUC vs max\_depth. You have to plot two curves, one curve from training data and another from cross-validation data in same plot. So that you will get a clear idea of when a model is overfitting or underfitting. You have to choose the hyperparameter before it overfits or underfits. Same can be followed for AUC vs min\_sample\_split.(Please sort the AUC values and plot if graph you are getting is varying too much) OR you can plot heatmaps or 3D-plots. you can find code from 3d scatter plot.ipynb**

meanly, if you can plot heatmap of 2D plots, you can find code from `eda_output_prompt.py` (Please upload the screen shots of 3D plots).

In [122]:

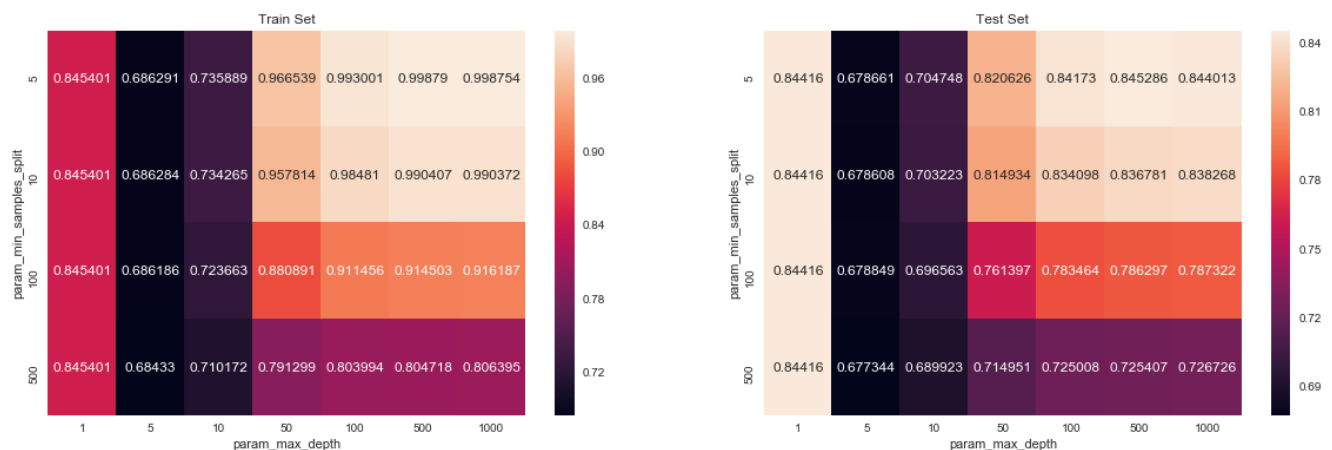
```
# https://seaborn.pydata.org/generated/seaborn.heatmap.html
# https://pandas.pydata.org/pandas-docs/stable/reference/api/pandas.DataFrame.unstack.html
import seaborn as sns; sns.set()
max_scores2=pd.DataFrame(modeltfidf.cv_results_).groupby(['param_min_samples_split','param_max_depth']).max().unstack(['mean_test_score','mean_train_score'])

fig,ax=plt.subplots(1,2,figsize=(20,6))

sns.heatmap(max_scores2.mean_train_score,annot=True,fmt='4g',ax=ax[0])
sns.heatmap(max_scores2.mean_test_score,annot=True,fmt='4g',ax=ax[1])

ax[0].set_title('Train Set')
ax[1].set_title('Test Set')

plt.show()
```



## Conclusion

For all the various values of `min_samples_split=5` and `max_depth=500` is giving the best score for test data.

In [188]:

```
#code source:
http://occam.olin.edu/sites/default/files/DataScienceMaterials/machine_learning_lecture_2/Machine%20Learning%20Lecture%202.html
from sklearn.model_selection import train_test_split
#from sklearn.grid_search import GridSearchCV
from sklearn.model_selection import GridSearchCV
from sklearn.datasets import *
from sklearn.tree import DecisionTreeClassifier

d_range=[500]
split_range=[5]

param_grid=dict(max_depth=d_range,min_samples_split=split_range)
print(param_grid)

#Using GridSearchCV
# https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.GridSearchCV.html
modeltfidfB = GridSearchCV(DecisionTreeClassifier(class_weight="balanced"), param_grid, scoring = 'f1', cv=5)
modeltfidfB.fit(X_tr_tfidf, y_train)

print(modeltfidfB.best_estimator_)
print(modeltfidfB.score(X_te_tfidf, y_test))

{'max_depth': [500], 'min_samples_split': [5]}
DecisionTreeClassifier(class_weight='balanced', criterion='gini',
                        max_depth=500, max_features=None, max_leaf_nodes=None,
                        min_impurity_decrease=0.0, min_impurity_split=None,
                        min_samples_leaf=1, min_samples_split=5,
                        min_weight_fraction_leaf=0.0, presort=False, random_state=None)
```

```

min_weight_fraction_leaf=0.0, presort=False, random_state=None,
splitter='best')
0.8456935630099729

```

In [190]:

```
best_tuned_parameters = [{'max_depth': [500], 'min_samples_split' : [5]}]
```

In [191]:

```

# https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve
from sklearn.metrics import roc_curve, auc

#model = GridSearchCV(LogisticRegression(), best_tuned_parameters)
modeltfidfB = GridSearchCV(DecisionTreeClassifier(), best_tuned_parameters)
modeltfidfB.fit(X_tr_tfidf, y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive
class
# not the predicted outputs
#print(type(model.predict_proba(X_tr_bow)))
#print(model.predict_proba(X_tr_bow))
#print(model.predict_proba(X_tr_bow)[: ,1])

y_train_tf_pred = modeltfidfB.predict_proba(X_tr_tfidf)[: ,1]
y_test_tf_pred = modeltfidfB.predict_proba(X_te_tfidf)[: ,1]

print(modeltfidfB.best_estimator_)
print(modeltfidfB.score(X_te_tfidf, y_test))

train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_tf_pred)
test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_tf_pred)

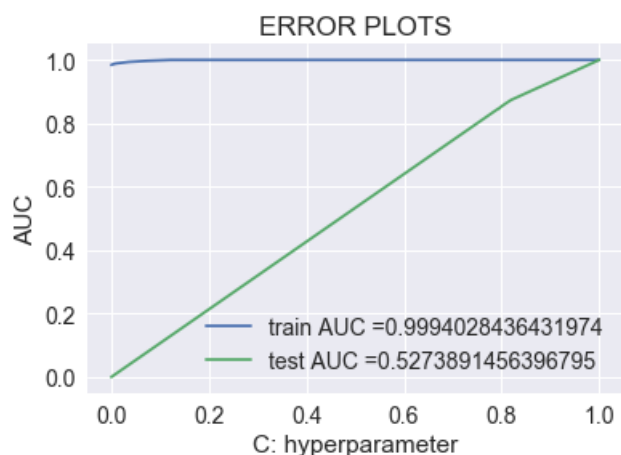
plt.plot(train_fpr, train_tpr, label="train AUC =" +str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="test AUC =" +str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel("C: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid(True)
plt.show()

```

```

DecisionTreeClassifier(class_weight=None, criterion='gini', max_depth=500,
max_features=None, max_leaf_nodes=None,
min_impurity_decrease=0.0, min_impurity_split=None,
min_samples_leaf=1, min_samples_split=5,
min_weight_fraction_leaf=0.0, presort=False, random_state=None,
splitter='best')
0.7492121212121212

```



In [192]:

```

# https://scikit-learn.org/stable/modules/generated/sklearn.metrics.confusion_matrix.html
print("="*100)
from sklearn.metrics import confusion_matrix
print("Train confusion matrix")

```

```
print('Train confusion matrix')
print(confusion_matrix(y_train, predict(y_train_tf_pred, tr_thresholds, train_fpr, train_tpr)))
print("Test confusion matrix")
print(confusion_matrix(y_test, predict(y_test_tf_pred, te_thresholds, test_fpr, test_tpr)))
```

```
=====

Train confusion matrix
the maximum value of tpr*(1-fpr) 0.9842580827333051 for threshold 1.0
[[ 5168    0]
 [ 446 27886]]
Test confusion matrix
the maximum value of tpr*(1-fpr) 0.17637054800527965 for threshold 1.0
[[ 532  2014]
 [ 2176 11778]]
```

In [193]:

```
import seaborn as snTr
import seaborn as snTe
import pandas as pdH
import matplotlib.pyplot as pltTr
import matplotlib.pyplot as pltTe

# https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
arrayTr=confusion_matrix(y_train, predict(y_train_tf_pred, tr_thresholds, train_fpr, train_tpr))
df_cmTr = pdH.DataFrame(arrayTr,range(2),range(2))
#print(arrayTr)
# https://stackoverflow.com/questions/32723798/how-do-i-add-a-title-to-seaborn-heatmap
axTr = pltTr.axes()

snTr.set(font_scale=1.4)#for label size

# https://seaborn.pydata.org/generated/seaborn.heatmap.html

snTr.heatmap(df_cmTr, annot=True,annot_kws={"size": 12},fmt="d",ax=axTr)# font size, format in
digit

labels=['Not Approved','Approved']
axTr.set_xticklabels(labels)
axTr.set_yticklabels(labels)
#Suggestion 4.Label confusion matrix heatmap with actual and predicted labels.
pltTr.title("Train confusion matrix")
pltTr.xlabel("Predicted")
pltTr.ylabel("Actual")
pltTr.show()

# https://stackoverflow.com/questions/50947776/plot-two-seaborn-heatmap-graphs-side-by-side
#fig, ax =plt.subplots(1,1)

# https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
arrayTe=confusion_matrix(y_test, predict(y_test_tf_pred, te_thresholds, test_fpr, test_tpr))
df_cmTe = pdH.DataFrame(arrayTe,range(2),range(2))

axTe = pltTe.axes()

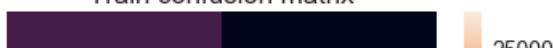
snTe.set(font_scale=1.4)#for label size

# https://seaborn.pydata.org/generated/seaborn.heatmap.html
snTe.heatmap(df_cmTe, annot=True,annot_kws={"size": 12},fmt="d",ax=axTe)# font size, format in
digit

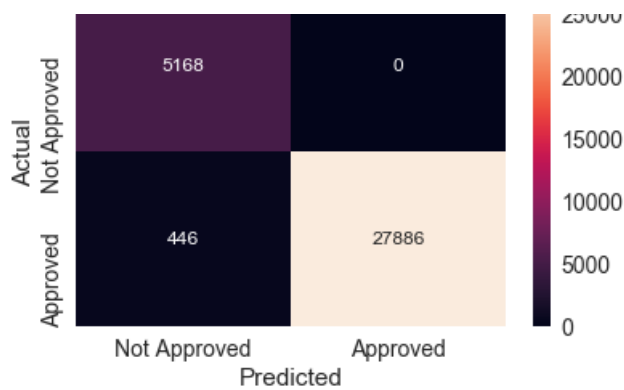
#Suggestion 4.Label confusion matrix heatmap with actual and predicted labels.
axTe.set_xticklabels(labels)
axTe.set_yticklabels(labels)
pltTe.title("Test confusion matrix")
pltTe.xlabel("Predicted")
pltTe.ylabel("Actual")
pltTe.show()
```

the maximum value of tpr\*(1-fpr) 0.9842580827333051 for threshold 1.0

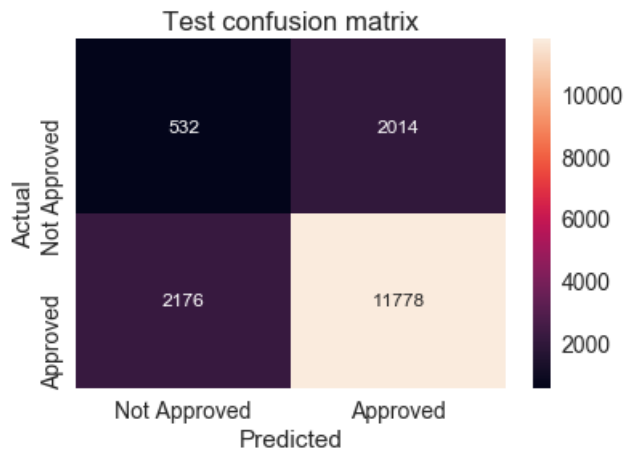
Train confusion matrix







the maximum value of  $\text{tpr} \times (1 - \text{fpr})$  0.17637054800527965 for threshold 1.0



1. Once after you plot the confusion matrix with the test data, get all the `false positive data points`

- Plot the WordCloud [WordCloud](#)
- Plot the box plot with the `price` of these `false positive data points`
- Plot the pdf with the `teacher\_number\_of\_previously\_posted\_projects` of these `false positive data points`

In [194]:

```
predict_tfidf = predict(y_test_tf_pred, te_thresholds, test_fpr, test_tpr) # <==

#converting predict_bow list into numpy array.
# https://likegeeks.com/numpy-array-tutorial/
import numpy as np
y_predict_tfidf = np.array(predict_tfidf)
#print(y_predict_tfidf)
#print(type(y_predict_tfidf))

# traversing y_test and y_predict, and created new list y_FP, which is 1 for FALSE_POSITIVE, and 0
for rest
y_tfidf_FP=[]
count_tfidf_FP=0
# https://www.geeksforgeeks.org/numpy-iterating-over-array/
for ac,pre in np.nditer([y_test,y_predict_tfidf]):
    if(ac==0 and pre==1):
        y_tfidf_FP.append(1)
        count_tfidf_FP+=1
    else:
        y_tfidf_FP.append(0)

print("False Positive count:",count_tfidf_FP)
##print(type(y_test))
##print(y_test)
##print(type(y_predict_tfidf))

#print(y_predict)
#print(type(y_FP))
#print(y_FP)
#converting predict_bow list into numpy array.
y_tfidf_FP=np.array(y_tfidf_FP)
print(type(y_tfidf_FP))
```

```

# print(y_test)
print(y_tfidf_FP)

```

```

the maximum value of tpr*(1-fpr) 0.17637054800527965 for threshold 1.0
False Positive count: 2014
<class 'numpy.ndarray'>
[0 0 0 ... 1 1 0]

```

In [195]:

```

#X_test.columns
#X_test#_working.columns

```

In [196]:

```

X_test_tfidf_working = X_test.copy(deep=True)  #<<==

#print("X_test_working length:",len(X_test_working))
#print(type(X_test))
#print(X_test.shape)
#print(type(X_test_tfidf_working))
#print(X_test_tfidf_working#.head(2))

# Adding 3 numpy array into dataframe
#print(type(y_test))
#print(type(y_predict_tfidf))
#print(type(y_tfidf_FP))
#print(y_test)
#print(y_predict_tfidf)
#print(y_tfidf_FP)
#
#print(type(X_test_tfidf_working))
X_test_tfidf_working['Y_Actual'] = y_test
print(X_test_tfidf_working.head(2))

```

```

      Unnamed: 0      id      teacher_id school_state \
33081      141937  p034157  0d7b3cd172c5b19f83a0ed303f46b729      AR
26184      33737  p225681  d44f8cb33de45fce2126bb699d302808      TN

      project_submitted_datetime \
33081      2016-09-26 16:20:26
26184      2016-12-26 14:31:27

      project_title \
33081  Scratching the Surface in Collabortive Learning!
26184      Ready 2 Read!

      project_essay_1 \
33081  With the upcoming school year, I am fortunate ...
26184  I am truly blessed to have the opportunity to ...

      project_essay_2 project_essay_3 \
33081  As my students become more technologically inc...      NaN
26184  My classroom consists of students on all learn...      NaN

      project_essay_4      ...      compound      price      quantity      essay_wc      essay_len \
33081      NaN      ...      0.9944      438.89      2      213      1528
26184      NaN      ...      0.9865      232.49      6      107      763

      title_wc      title_len      prj_res_sum_wc      prj_res_sum_len      Y_Actual
33081      4      40      14      106      1
26184      3      12      9      61      1

```

[2 rows x 31 columns]

In [197]:

```

# how to add numpy array into dataframe | https://stackoverflow.com/questions/26666919/add-column-
in-dataframe-from-list
print(X_test_tfidf_working.columns)
X_test_tfidf_working['y_predict_tfidf'] = y_predict_tfidf
X_test_tfidf_working['y_tfidf_FP'] = y_tfidf_FP
print(X_test_tfidf_working.columns)

```

```
Index(['Unnamed: 0', 'id', 'teacher_id', 'school_state',
      'project_submitted_datetime', 'project_title', 'project_essay_1',
      'project_essay_2', 'project_essay_3', 'project_essay_4',
      'project_resource_summary',
      'teacher_number_of_previously_posted_projects', 'project_is_approved',
      'clean_categories', 'clean_subcategories', 'clean_grade',
      'clean_tea_pfx', 'essay', 'neg', 'pos', 'neu', 'compound', 'price',
      'quantity', 'essay_wc', 'essay_len', 'title_wc', 'title_len',
      'prj_res_sum_wc', 'prj_res_sum_len', 'Y_Actual'],
      dtype='object')
Index(['Unnamed: 0', 'id', 'teacher_id', 'school_state',
      'project_submitted_datetime', 'project_title', 'project_essay_1',
      'project_essay_2', 'project_essay_3', 'project_essay_4',
      'project_resource_summary',
      'teacher_number_of_previously_posted_projects', 'project_is_approved',
      'clean_categories', 'clean_subcategories', 'clean_grade',
      'clean_tea_pfx', 'essay', 'neg', 'pos', 'neu', 'compound', 'price',
      'quantity', 'essay_wc', 'essay_len', 'title_wc', 'title_len',
      'prj_res_sum_wc', 'prj_res_sum_len', 'Y_Actual', 'y_predict_tfidf',
      'y_tfidf_FP'],
      dtype='object')
```

In [198]:

```
X_test_tfidf_FP=X_test_tfidf_working.copy(deep=True)
X_test_tfidf_FP.teacher_number_of_previously_posted_projects
X_test_tfidf_FP.price
X_test_tfidf_FP.columns
```

Out[198]:

```
Index(['Unnamed: 0', 'id', 'teacher_id', 'school_state',
      'project_submitted_datetime', 'project_title', 'project_essay_1',
      'project_essay_2', 'project_essay_3', 'project_essay_4',
      'project_resource_summary',
      'teacher_number_of_previously_posted_projects', 'project_is_approved',
      'clean_categories', 'clean_subcategories', 'clean_grade',
      'clean_tea_pfx', 'essay', 'neg', 'pos', 'neu', 'compound', 'price',
      'quantity', 'essay_wc', 'essay_len', 'title_wc', 'title_len',
      'prj_res_sum_wc', 'prj_res_sum_len', 'Y_Actual', 'y_predict_tfidf',
      'y_tfidf_FP'],
      dtype='object')
```

In [199]:

```
X_test_tfidf_FP=X_test_tfidf_FP.loc[:,["teacher_number_of_previously_posted_projects","essay","pri
ce","y_tfidf_FP"]]
X_test_tfidf_FP.head(2)
```

Out[199]:

	teacher_number_of_previously_posted_projects	essay	price	y_tfidf_FP
33081	43	With the upcoming school year, I am fortunate ...	438.89	0
26184	1	I am truly blessed to have the opportunity to ...	232.49	0

In [200]:

```
X_test_tfidf_FP = X_test_tfidf_FP[X_test_tfidf_FP.y_tfidf_FP==True]
```

In [201]:

```
# https://www.geeksforgeeks.org/generating-word-cloud-python/
# Python program to generate WordCloud

# importing all necessary modules
from wordcloud import WordCloud, STOPWORDS
import matplotlib.pyplot as plt
import pandas as pd
```

```
# Reads 'Youtube04-Eminem.csv' file
#df = pd.read_csv(r"Youtube04-Eminem.csv", encoding ="latin-1")

comment_words = ' '
stopwords = set(STOPWORDS)

# iterate through the csv file
for val in X_test_tfidf_FP["essay"][:1]:

    # typecaste each val to string
    val = str(val)

    # split the value
    tokens = val.split()

    # Converts each token into lowercase
    for i in range(len(tokens)):
        tokens[i] = tokens[i].lower()

    for words in tokens:
        comment_words = comment_words + words + ' '

wordcloud = WordCloud(width = 800, height = 800,
                        background_color='white',
                        stopwords = stopwords,
                        min_font_size = 10).generate(comment_words)

# plot the WordCloud image
plt.figure(figsize = (8, 8), facecolor = None)
plt.imshow(wordcloud)
plt.axis("off")
plt.tight_layout(pad = 0)

plt.show()
```

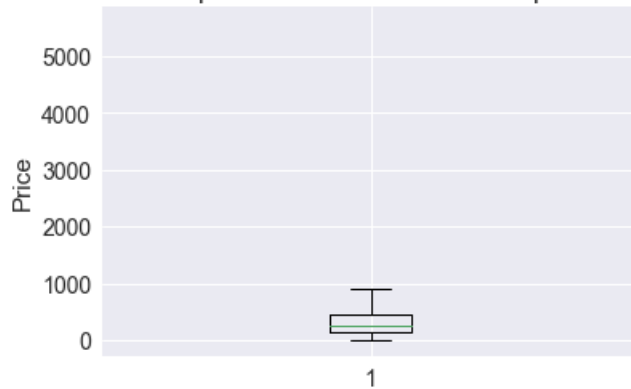


**Plot the box plot with the price of these false positive data points**

In [202]:

```
import matplotlib.pyplot as plt
# https://glowingpython.blogspot.com/2012/09/boxplot-with-matplotlib.html
plt.boxplot([X_test_tfidf_FP["price"]])
plt.title('Box Plots with the price of these TFIDF\'s false positive data points')
#labels = ('Price')
#plt.xticks([1],labels,rotation=90)
plt.ylabel('Price')
plt.grid(True)
plt.show()
```

Box Plots with the price of these TFIDF's false positive data points



**Plot the pdf with the teacher\_number\_of\_previously\_posted\_projects of these false positive data points**

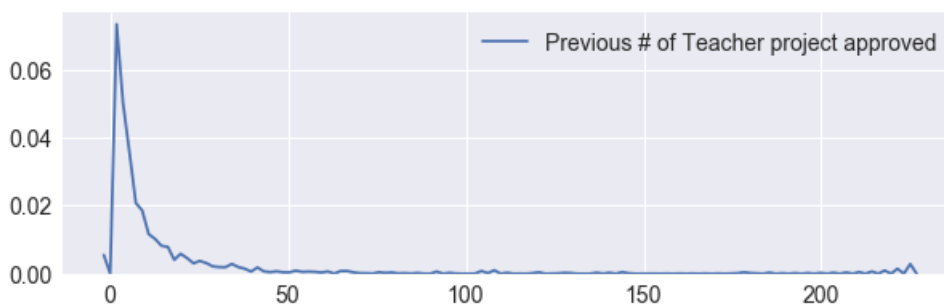
In [203]:

```
plt.figure(figsize=(10,3))

# https://seaborn.pydata.org/generated/seaborn.kdeplot.html | kernel density estimate | sns.kdeplot
# bw : {'scott' | 'silverman' | scalar | pair of scalars }, optional
# Name of reference method to determine kernel size, scalar factor, or scalar for each dimension
# of the bivariate plot. Note that the underlying computational libraries have different
# interpretations
# for this parameter: statsmodels uses it directly, but scipy treats it as a scaling factor
# for the standard deviation of the data.

sns.kdeplot(X_test_tfidf_FP["teacher_number_of_previously_posted_projects"],label="Previous # of T
eacher project approved", bw=0.6)
plt.legend()
plt.show()

#PDF
```



#### 2.4.2.1 Graphviz visualization of Decision Tree on TFIDF, SET 2

In [204]:

```
#print(len(aa))
#print(len(b))
#print(len(c))
```



```

        feature_names=feature_list,
        class_names=True,
        filled=True, rounded=True,
        special_characters=False)

graph = g.Source(dot_data2)
graph

# https://scikit-learn.org/stable/modules/generated/sklearn.tree.export_graphviz.html#sklearn.tree.export_graphviz
# https://stackoverflow.com/questions/5316206/convert-dot-to-png-in-python

(graph,) = pydot.graph_from_dot_file('X_tr_tfidf_graphviz.dot')
graph.write_png('X_tr_tfidf_graphviz.png')

```

## 2.4.3 Applying Decision Trees on AVG W2V, SET 3

### 1. Apply Decision Tree Classifier(DecisionTreeClassifier) on these feature sets

- **Set 3:** categorical, numerical features + project\_title(AVG W2V)+ preprocessed\_eassay (AVG W2V)

In [142]:

```

# Please write all the code with proper documentation
# merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
from scipy.sparse import hstack
X_tr_avgW2V = hstack((X_train_essay_avg_w2v, X_train_title_avg_w2v, X_train_state_ohe,
X_train_clean_ohe, X_train_cleanSub_ohe, X_train_grade_ohe, X_train_teacher_ohe,
X_train_prjResSum_ohe, X_train_quantity_norm, X_train_TprevPrj_norm, X_train_price_norm)).tocsr()
X_te_avgW2V = hstack((X_test_essay_avg_w2v, X_test_title_avg_w2v, X_test_state_ohe, X_test_clean_ohe,
X_test_cleanSub_ohe, X_test_grade_ohe, X_test_teacher_ohe, X_test_prjResSum_ohe, X_test_quantity_norm,
X_test_TprevPrj_norm, X_test_price_norm)).tocsr()

print("Final Data matrix | Avg W2V")
print(X_tr_avgW2V.shape, y_train.shape)
print(X_te_avgW2V.shape, y_test.shape)
print("="*100)

```

```

Final Data matrix | Avg W2V
(33500, 11315) (33500,)
(16500, 11315) (16500,)
=====

```

In [143]:

```

#code source:
http://occam.olin.edu/sites/default/files/DataScienceMaterials/machine_learning_lecture_2/Machine%20Learning%20Lecture%202.html
from sklearn.model_selection import train_test_split
from sklearn.grid_search import GridSearchCV
from sklearn.model_selection import GridSearchCV
from sklearn.datasets import *
from sklearn.tree import DecisionTreeClassifier

d_range=[1,5,10,50,100,500,1000]
split_range=[5,10,100,500]

param_grid=dict(max_depth=d_range,min_samples_split=split_range)
print(param_grid)

#Using GridSearchCV
# https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.GridSearchCV.html
modelavgW2V = GridSearchCV(DecisionTreeClassifier(class_weight="balanced"), param_grid, scoring =
'f1', cv=5)
modelavgW2V.fit(X_tr_avgW2V, y_train)

print(modelavgW2V.best_estimator_)
print(modelavgW2V.score(X_te_avgW2V, y_test))

```

```

{'max_depth': [1, 5, 10, 50, 100, 500, 1000], 'min_samples_split': [5, 10, 100, 500]}
DecisionTreeClassifier(class_weight='balanced', criterion='gini', max_depth=1,
max_features=None, max_leaf_nodes=None,

```

```

min_impurity_decrease=0.0, min_impurity_split=None,
min_samples_leaf=1, min_samples_split=5,
min_weight_fraction_leaf=0.0, presort=False, random_state=None,
splitter='best')
0.8976886471787899

```

## Conclusion

For all the various values of min\_samples\_split=5 and max\_depth=1 is giving the best score for test data.

In [144]:

```

#code source:
http://occam.olin.edu/sites/default/files/DataScienceMaterials/machine_learning_lecture_2/Machine%20Learning%20Lecture%202.html
from sklearn.model_selection import train_test_split
#from sklearn.grid_search import GridSearchCV
from sklearn.model_selection import GridSearchCV
from sklearn.datasets import *
from sklearn.tree import DecisionTreeClassifier

d_range=[1]
split_range=[5]

param_grid=dict(max_depth=d_range,min_samples_split=split_range)
print(param_grid)

#Using GridSearchCV
# https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.GridSearchCV.html
modelavgW2VB = GridSearchCV(DecisionTreeClassifier(class_weight="balanced"), param_grid, scoring =
'f1', cv=5)
modelavgW2VB.fit(X_tr_avgW2V, y_train)

print(modelavgW2VB.best_estimator_)
print(modelavgW2VB.score(X_te_avgW2V, y_test))

{'max_depth': [1], 'min_samples_split': [5]}
DecisionTreeClassifier(class_weight='balanced', criterion='gini', max_depth=1,
max_features=None, max_leaf_nodes=None,
min_impurity_decrease=0.0, min_impurity_split=None,
min_samples_leaf=1, min_samples_split=5,
min_weight_fraction_leaf=0.0, presort=False, random_state=None,
splitter='best')
0.8976886471787899

```

In [145]:

```
best_tuned_parameters = [{'max_depth': [1], 'min_samples_split' : [5]}]
```

In [146]:

```

# https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve
from sklearn.metrics import roc_curve, auc

#model = GridSearchCV(LogisticRegression(), best_tuned_parameters)
modelavgW2VB = GridSearchCV(DecisionTreeClassifier(), best_tuned_parameters)
modelavgW2VB.fit(X_tr_avgW2V, y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive
class
# not the predicted outputs
#print(type(model.predict_proba(X_tr_avgW2V)))
#print(model.predict_proba(X_tr_avgW2V))
#print(model.predict_proba(X_tr_avgW2V)[:,1])

y_train_avgW2V_pred = modelavgW2VB.predict_proba(X_tr_avgW2V)[:,1]
y_test_avgW2V_pred = modelavgW2VB.predict_proba(X_te_avgW2V)[:,1]

print(modelavgW2VB.best_estimator_)
print(modelavgW2VB.score(X_te_avgW2V, y_test))

train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_avgW2V_pred)
test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_avgW2V_pred)

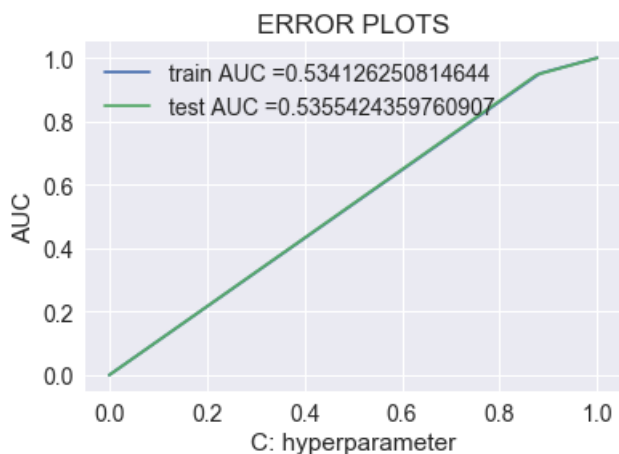
```



```
test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_avgwzv_pred)

plt.plot(train_fpr, train_tpr, label="train AUC =" + str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="test AUC =" + str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel("C: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid(True)
plt.show()
```

```
DecisionTreeClassifier(class_weight=None, criterion='gini', max_depth=1,
                        max_features=None, max_leaf_nodes=None,
                        min_impurity_decrease=0.0, min_impurity_split=None,
                        min_samples_leaf=1, min_samples_split=5,
                        min_weight_fraction_leaf=0.0, presort=False, random_state=None,
                        splitter='best')
0.8456969696969697
```



In [208]:

```
# https://scikit-learn.org/stable/modules/generated/sklearn.metrics.confusion_matrix.html
print("="*100)
from sklearn.metrics import confusion_matrix
print("Train confusion matrix")
print(confusion_matrix(y_train, predict(y_train_avgW2V_pred, tr_thresholds, train_fpr,
train_tpr)))
print("Test confusion matrix")
print(confusion_matrix(y_test, predict(y_test_avgW2V_pred, te_thresholds, test_fpr, test_tpr)))
```

```
=====
Train confusion matrix
the maximum value of tpr*(1-fpr) 0.9842580827333051 for threshold 1.0
[[ 5168    0]
 [28332    0]]
Test confusion matrix
the maximum value of tpr*(1-fpr) 0.17637054800527965 for threshold 1.0
[[ 2546    0]
 [13954    0]]
```

In [209]:

```
import seaborn as snTr
import seaborn as snTe
import pandas as pdH
import matplotlib.pyplot as pltTr
import matplotlib.pyplot as pltTe

# https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
arrayTr=confusion_matrix(y_train, predict(y_train_avgW2V_pred, tr_thresholds, train_fpr, train_tpr
))
df_cmTr = pdH.DataFrame(arrayTr,range(2),range(2))
#print(arrayTr)
# https://stackoverflow.com/questions/32723798/how-do-i-add-a-title-to-seaborn-heatmap
axTr = pltTr.axes()
```

```

snTr.set(font_scale=1.4)#for label size

# https://seaborn.pydata.org/generated/seaborn.heatmap.html

snTr.heatmap(df_cmTr, annot=True,annot_kws={"size": 12},fmt="d",ax=axTr)# font size, format in
digit

labels=['Not Approved','Approved']
axTr.set_xticklabels(labels)
axTr.set_yticklabels(labels)
#Suggestion 4.Label confusion matrix heatmap with actual and predicted labels.
pltTr.title("Train confusion matrix")
pltTr.xlabel("Predicted")
pltTr.ylabel("Actual")
pltTr.show()

# https://stackoverflow.com/questions/50947776/plot-two-seaborn-heatmap-graphs-side-by-side
#fig, ax =plt.subplots(1,1)

# https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
arrayTe=confusion_matrix(y_test, predict(y_test_avgW2V_pred, te_thresholds, test_fpr, test_tpr))
df_cmTe = pdH.DataFrame(arrayTe,range(2),range(2))

axTe = pltTe.axes()

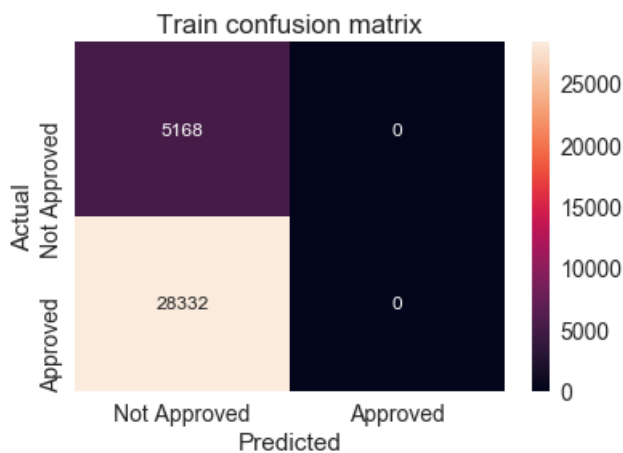
snTe.set(font_scale=1.4)#for label size

# https://seaborn.pydata.org/generated/seaborn.heatmap.html
snTe.heatmap(df_cmTe, annot=True,annot_kws={"size": 12},fmt="d",ax=axTe)# font size, format in
digit

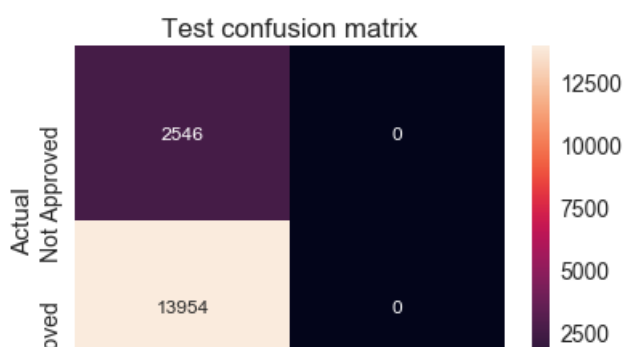
#Suggestion 4.Label confusion matrix heatmap with actual and predicted labels.
axTe.set_xticklabels(labels)
axTe.set_yticklabels(labels)
pltTe.title("Test confusion matrix")
pltTe.xlabel("Predicted")
pltTe.ylabel("Actual")
pltTe.show()

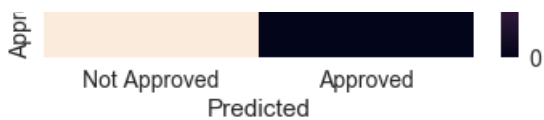
```

the maximum value of  $tpr \cdot (1 - fpr)$  0.9842580827333051 for threshold 1.0



the maximum value of  $tpr \cdot (1 - fpr)$  0.17637054800527965 for threshold 1.0





- Once after you plot the confusion matrix with the test data, get all the `false positive data points`
  - Plot the WordCloud [WordCloud](#)
  - Plot the box plot with the `price` of these `false positive data points`
  - Plot the pdf with the `teacher\_number\_of\_previously\_posted\_projects` of these `false positive data points`

In [210]:

```
predict_avgW2V = predict(y_test_avgW2V_pred, te_thresholds, test_fpr, test_tpr) # <==

#converting predict_bow list into numpy array.
# https://likegeeks.com/numpy-array-tutorial/
import numpy as np
y_predict_avgW2V = np.array(predict_avgW2V)
#print(y_predict_avgW2V)
#print(type(y_predict_avgW2V))

# traversing y_test and y_predict, and created new list y_FP, which is 1 for FALSE_POSITIVE, and 0
for rest
y_avgW2V_FP=[]
count_avgW2V_FP=0
# https://www.geeksforgeeks.org/numpy-iterating-over-array/
for ac,pre in np.nditer([y_test,y_predict_avgW2V]):
    if(ac==0 and pre==1):
        y_avgW2V_FP.append(1)
        count_avgW2V_FP+=1
    else:
        y_avgW2V_FP.append(0)

print("False Positive count:",count_avgW2V_FP)
##print(type(y_test))
##print(y_test)
##print(type(y_predict_avgW2V))

#print(y_predict)
#print(type(y_FP))
#print(y_FP)
#converting predict_bow list into numpy array.
y_avgW2V_FP=np.array(y_avgW2V_FP)
print(type(y_avgW2V_FP))
print(y_avgW2V_FP)
```

the maximum value of tpr\*(1-fpr) 0.17637054800527965 for threshold 1.0  
 False Positive count: 0  
 <class 'numpy.ndarray'>  
 [0 0 0 ... 0 0 0]

In [211]:

```
X_test_avgW2V_working = X_test.copy(deep=True) #<==

#print("X_test_working length:",len(X_test_working))
#print(type(X_test))
#print(X_test.shape)
#print(type(X_test_avgW2V_working))
#print(X_test_avgW2V_working#.head(2))

# Adding 3 numpy array into dataframe
#print(type(y_test))
#print(type(y_predict_avgW2V))
#print(type(y_avgW2V_FP))
#print(y_test)
#print(y_predict_avgW2V)
#print(y_avgW2V_FP)
#
#print(type(X_test_avgW2V_working))
X_test_avgW2V_working['Y_Actual'] = y_test
print(X_test_avgW2V_working.head(2))
```

	Unnamed: 0	id	teacher_id	school_state	
33081	141937	p034157	0d7b3cd172c5b19f83a0ed303f46b729	AR	
26184	33737	p225681	d44f8cb33de45fce2126bb699d302808	TN	

	project_submitted_datetime	
33081	2016-09-26 16:20:26	
26184	2016-12-26 14:31:27	

	project_title	
33081	Scratching the Surface in Collaborative Learning!	
26184	Ready 2 Read!	

	project_essay_1	
33081	With the upcoming school year, I am fortunate ...	
26184	I am truly blessed to have the opportunity to ...	

	project_essay_2	project_essay_3	
33081	As my students become more technologically inc...	NaN	
26184	My classroom consists of students on all learn...	NaN	

	project_essay_4	...	compound	price	quantity	essay_wc	essay_len	
33081	NaN	...	0.9944	438.89	2	213	1528	
26184	NaN	...	0.9865	232.49	6	107	763	

	title_wc	title_len	prj_res_sum_wc	prj_res_sum_len	Y_Actual
33081	4	40	14	106	1
26184	3	12	9	61	1

[2 rows x 31 columns]

In [212]:

```
# how to add numpy array into dataframe | https://stackoverflow.com/questions/26666919/add-column-in-dataframe-from-list
print(X_test_avgW2V_working.columns)
X_test_avgW2V_working['y_predict_avgW2V'] = y_predict_avgW2V
X_test_avgW2V_working['y_avgW2V_FP'] = y_avgW2V_FP
print(X_test_avgW2V_working.columns)
```

```
Index(['Unnamed: 0', 'id', 'teacher_id', 'school_state',
      'project_submitted_datetime', 'project_title', 'project_essay_1',
      'project_essay_2', 'project_essay_3', 'project_essay_4',
      'project_resource_summary',
      'teacher_number_of_previously_posted_projects', 'project_is_approved',
      'clean_categories', 'clean_subcategories', 'clean_grade',
      'clean_tea_pfx', 'essay', 'neg', 'pos', 'neu', 'compound', 'price',
      'quantity', 'essay_wc', 'essay_len', 'title_wc', 'title_len',
      'prj_res_sum_wc', 'prj_res_sum_len', 'Y_Actual'],
      dtype='object')
Index(['Unnamed: 0', 'id', 'teacher_id', 'school_state',
      'project_submitted_datetime', 'project_title', 'project_essay_1',
      'project_essay_2', 'project_essay_3', 'project_essay_4',
      'project_resource_summary',
      'teacher_number_of_previously_posted_projects', 'project_is_approved',
      'clean_categories', 'clean_subcategories', 'clean_grade',
      'clean_tea_pfx', 'essay', 'neg', 'pos', 'neu', 'compound', 'price',
      'quantity', 'essay_wc', 'essay_len', 'title_wc', 'title_len',
      'prj_res_sum_wc', 'prj_res_sum_len', 'Y_Actual', 'y_predict_avgW2V',
      'y_avgW2V_FP'],
      dtype='object')
```

In [213]:

```
X_test_avgW2V_FP=X_test_avgW2V_working.copy(deep=True)
X_test_avgW2V_FP.teacher_number_of_previously_posted_projects
X_test_avgW2V_FP.price
X_test_avgW2V_FP.columns
```

Out [213]:

```
Index(['Unnamed: 0', 'id', 'teacher_id', 'school_state',
      'project_submitted_datetime', 'project_title', 'project_essay_1',
      'project_essay_2', 'project_essay_3', 'project_essay_4',
      'project_resource_summary',
```

```

project_resource_summary,
'teacher_number_of_previously_posted_projects', 'project_is_approved',
'clean_categories', 'clean_subcategories', 'clean_grade',
'clean_tea_pfx', 'essay', 'neg', 'pos', 'neu', 'compound', 'price',
'quantity', 'essay_wc', 'essay_len', 'title_wc', 'title_len',
'prj_res_sum_wc', 'prj_res_sum_len', 'Y_Actual', 'y_predict_avgW2V',
'y_avgW2V_FP'],
dtype='object')

```

In [214]:

```

X_test_avgW2V_FP=X_test_avgW2V_FP.loc[:,["teacher_number_of_previously_posted_projects","essay","p
rice","y_avgW2V_FP"]]
print(X_test_avgW2V_FP.head(2))

X_test_avgW2V_FP = X_test_avgW2V_FP[X_test_avgW2V_FP.y_avgW2V_FP==True]

```

```

teacher_number_of_previously_posted_projects  \
33081                                         43
26184                                         1

essay  price  y_avgW2V_FP
33081  With the upcoming school year, I am fortunate ...  438.89      0
26184  I am truly blessed to have the opportunity to ...  232.49      0

```

In [154]:

```

# https://www.geeksforgeeks.org/generating-word-cloud-python/
# Python program to generate WordCloud

# importing all necessary modules
from wordcloud import WordCloud, STOPWORDS
import matplotlib.pyplot as plt
import pandas as pd

# Reads 'Youtube04-Eminem.csv' file
df = pd.read_csv(r"Youtube04-Eminem.csv", encoding = "latin-1")

comment_words = ' '
stopwords = set(STOPWORDS)

# iterate through the csv file
for val in X_test_avgW2V_FP["essay"][:1]:

    # typecaste each val to string
    val = str(val)

    # split the value
    tokens = val.split()

    # Converts each token into lowercase
    for i in range(len(tokens)):
        tokens[i] = tokens[i].lower()

    for words in tokens:
        comment_words = comment_words + words + ' '

wordcloud = WordCloud(width = 800, height = 800,
                      background_color = 'white',
                      stopwords = stopwords,
                      min_font_size = 10).generate(comment_words)

# plot the WordCloud image
plt.figure(figsize = (8, 8), facecolor = None)
plt.imshow(wordcloud)
plt.axis("off")
plt.tight_layout(pad = 0)

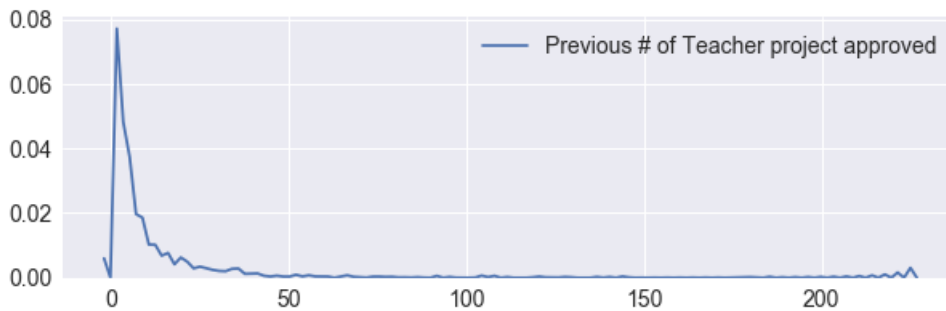
plt.show()

```



```
t
# bw : {'scott' | 'silverman' | scalar | pair of scalars }, optional
# Name of reference method to determine kernel size, scalar factor, or scalar for each dimension
# of the bivariate plot. Note that the underlying computational libraries have different
# interpretations
# for this parameter: statsmodels uses it directly, but scipy treats it as a scaling factor
# for the standard deviation of the data.

sns.kdeplot(X_test_avgW2V_FP["teacher_number_of_previously_posted_projects"],label="Previous # of
Teacher project approved", bw=0.6)
plt.legend()
plt.show()
```



## 2.4.4 Applying Decision Trees on TFIDF W2V, SET 4

### 1. Apply Decision Tree Classifier(DecisionTreeClassifier) on these feature sets

- **Set 4:** categorical, numerical features + project\_title(TFIDF W2V)+ preprocessed\_eassay (TFIDF W2V)

In [157]:

```
# merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
from scipy.sparse import hstack
X_tr_tfidf_W2V = hstack((tr_tfidf_w2v_essay_vectors, tr_tfidf_w2v_title_vectors, X_train_state_ohe,
X_train_clean_ohe, X_train_cleanSub_ohe, X_train_grade_ohe, X_train_teacher_ohe,
X_train_prjResSum_ohe, X_train_quantity_norm, X_train_TprevPrj_norm, X_train_price_norm)).tocsr()
X_te_tfidf_W2V = hstack((te_tfidf_w2v_essay_vectors, te_tfidf_w2v_title_vectors, X_test_state_ohe,
X_test_clean_ohe, X_test_cleanSub_ohe, X_test_grade_ohe, X_test_teacher_ohe, X_test_prjResSum_ohe,
X_test_quantity_norm, X_test_TprevPrj_norm, X_test_price_norm)).tocsr()

print("Final Data matrix | TFIDF W2V")
print(X_tr_tfidf_W2V.shape, y_train.shape)
print(X_te_tfidf_W2V.shape, y_test.shape)
print("="*100)
```

```
Final Data matrix | TFIDF W2V
(33500, 11315) (33500,)
(16500, 11315) (16500,)
```

In [158]:

```
#code source:
http://occam.olin.edu/sites/default/files/DataScienceMaterials/machine_learning_lecture_2/Machine%
rning%20Lecture%202.html
from sklearn.model_selection import train_test_split
#from sklearn.grid_search import GridSearchCV
from sklearn.model_selection import GridSearchCV
from sklearn.datasets import *
from sklearn.tree import DecisionTreeClassifier

d_range=[1,5,10,50,100,500,1000]
split_range=[5,10,100,500]

param_grid=dict(max_depth=d_range,min_samples_split=split_range)
print(param_grid)

#Using GridSearchCV
```

```
# https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.GridSearchCV.html
modeltfidf_W2V = GridSearchCV(DecisionTreeClassifier(class_weight="balanced"), param_grid, scoring
= 'f1', cv=5)
modeltfidf_W2V.fit(X_tr_tfidf_W2V, y_train)

print(modeltfidf_W2V.best_estimator_)
print(modeltfidf_W2V.score(X_te_tfidf_W2V, y_test))
```

```
{'max_depth': [1, 5, 10, 50, 100, 500, 1000], 'min_samples_split': [5, 10, 100, 500]}
DecisionTreeClassifier(class_weight='balanced', criterion='gini', max_depth=1,
    max_features=None, max_leaf_nodes=None,
    min_impurity_decrease=0.0, min_impurity_split=None,
    min_samples_leaf=1, min_samples_split=5,
    min_weight_fraction_leaf=0.0, presort=False, random_state=None,
    splitter='best')
0.8976886471787899
```

## Conclusion

For all the various values of min\_samples\_split=5 and max\_depth=1 is giving the best score for test data.

In [159]:

```
#code source:
http://occam.olin.edu/sites/default/files/DataScienceMaterials/machine_learning_lecture_2/Machine%20Learning%20Lecture%202.html
from sklearn.model_selection import train_test_split
from sklearn.grid_search import GridSearchCV
from sklearn.model_selection import GridSearchCV
from sklearn.datasets import *
from sklearn.tree import DecisionTreeClassifier

d_range=[1]
split_range=[5]

param_grid=dict(max_depth=d_range,min_samples_split=split_range)
print(param_grid)

#Using GridSearchCV
# https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.GridSearchCV.html
modeltfidf_W2VB = GridSearchCV(DecisionTreeClassifier(class_weight="balanced"), param_grid,
scoring = 'f1', cv=5)
modeltfidf_W2VB.fit(X_tr_tfidf_W2V, y_train)

print(modeltfidf_W2VB.best_estimator_)
print(modeltfidf_W2VB.score(X_te_tfidf_W2V, y_test))
```

```
{'max_depth': [1], 'min_samples_split': [5]}
DecisionTreeClassifier(class_weight='balanced', criterion='gini', max_depth=1,
    max_features=None, max_leaf_nodes=None,
    min_impurity_decrease=0.0, min_impurity_split=None,
    min_samples_leaf=1, min_samples_split=5,
    min_weight_fraction_leaf=0.0, presort=False, random_state=None,
    splitter='best')
0.8976886471787899
```

In [160]:

```
best_tuned_parameters = [{'max_depth': [1], 'min_samples_split' :[5]}]
```

In [161]:

```
# https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve
from sklearn.metrics import roc_curve, auc

#model = GridSearchCV(LogisticRegression(), best_tuned_parameters)
modeltfidf_W2VB = GridSearchCV(DecisionTreeClassifier(), best_tuned_parameters)
modeltfidf_W2VB.fit(X_tr_tfidf_W2V, y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive
class
# not the predicted outputs
```



```

# Now we produce the output
#print(type(model.predict_proba(X_tr_bow)))
#print(model.predict_proba(X_tr_bow))
#print(model.predict_proba(X_tr_bow)[:,:1])

y_train_tfidf_w2v_pred = modeltfidf_W2VB.predict_proba(X_tr_tfidf_W2V)[:,:1]
y_test_tfidf_w2v_pred = modeltfidf_W2VB.predict_proba(X_te_tfidf_W2V)[:,:1]

print(modeltfidf_W2VB.best_estimator_)
print(modeltfidf_W2VB.score(X_te_tfidf_W2V, y_test))

train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_tfidf_w2v_pred)
test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_tfidf_w2v_pred)

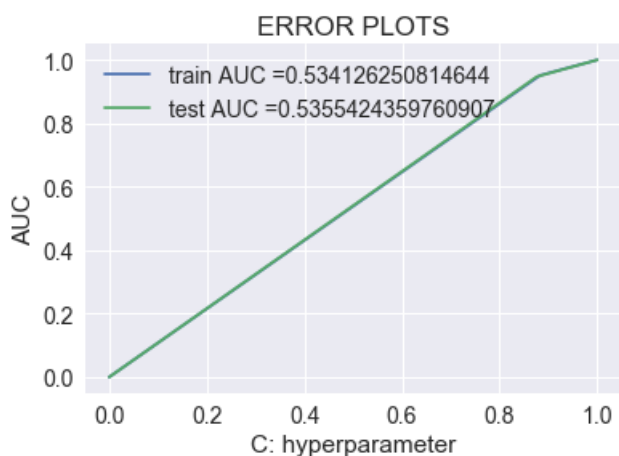
plt.plot(train_fpr, train_tpr, label="train AUC =" + str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="test AUC =" + str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel("C: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid(True)
plt.show()

```

```

DecisionTreeClassifier(class_weight=None, criterion='gini', max_depth=1,
                        max_features=None, max_leaf_nodes=None,
                        min_impurity_decrease=0.0, min_impurity_split=None,
                        min_samples_leaf=1, min_samples_split=5,
                        min_weight_fraction_leaf=0.0, presort=False, random_state=None,
                        splitter='best')
0.8456969696969697

```



In [162]:

```

# https://scikit-learn.org/stable/modules/generated/sklearn.metrics.confusion_matrix.html
print("="*100)
from sklearn.metrics import confusion_matrix
print("Train confusion matrix")
print(confusion_matrix(y_train, predict(y_train_tfidf_w2v_pred, tr_thresholds, train_fpr,
train_tpr)))
print("Test confusion matrix")
print(confusion_matrix(y_test, predict(y_test_tfidf_w2v_pred, te_thresholds, test_fpr, test_tpr)))

```

```

=====
Train confusion matrix
the maximum value of tpr*(1-fpr) 0.11135206899920401 for threshold 0.855
[[ 605  4563]
 [ 1383 26949]]
Test confusion matrix
the maximum value of tpr*(1-fpr) 0.1175342031122122 for threshold 0.855
[[ 316  2230]
 [ 740 13214]]

```

In [163]:

```

import seaborn as snTr

```

```

import seaborn as snTe
import pandas as pdH
import matplotlib.pyplot as pltTr
import matplotlib.pyplot as pltTe

# https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
arrayTr=confusion_matrix(y_train, predict(y_train_tfidf_w2v_pred, tr_thresholds, train_fpr,
train_tpr))
df_cmTr = pdH.DataFrame(arrayTr,range(2),range(2))
#print(arrayTr)
# https://stackoverflow.com/questions/32723798/how-do-i-add-a-title-to-seaborn-heatmap
axTr = pltTr.axes()

snTr.set(font_scale=1.4)#for label size

# https://seaborn.pydata.org/generated/seaborn.heatmap.html

snTr.heatmap(df_cmTr, annot=True,annot_kws={"size": 12},fmt="d",ax=axTr)# font size, format in
digit

labels=['Not Approved','Approved']
axTr.set_xticklabels(labels)
axTr.set_yticklabels(labels)
#Suggestion 4.Label confusion matrix heatmap with actual and predicted labels.
pltTr.title("Train confusion matrix")
pltTr.xlabel("Predicted")
pltTr.ylabel("Actual")
pltTr.show()

# https://stackoverflow.com/questions/50947776/plot-two-seaborn-heatmap-graphs-side-by-side
#fig, ax =plt.subplots(1,1)

# https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
arrayTe=confusion_matrix(y_test, predict(y_test_tfidf_w2v_pred, te_thresholds, test_fpr, test_tpr))
df_cmTe = pdH.DataFrame(arrayTe,range(2),range(2))

axTe = pltTe.axes()

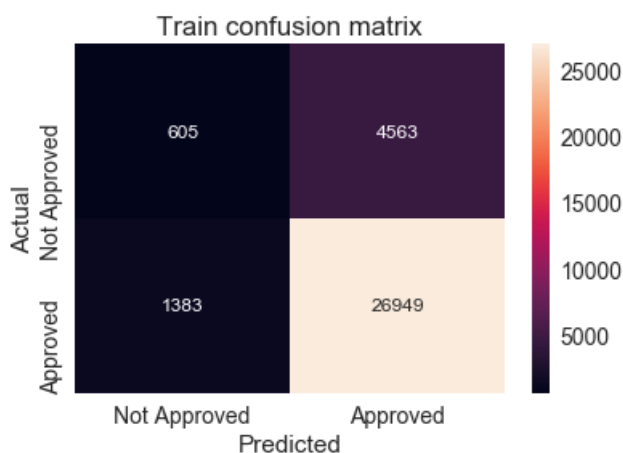
snTe.set(font_scale=1.4)#for label size

# https://seaborn.pydata.org/generated/seaborn.heatmap.html
snTe.heatmap(df_cmTe, annot=True,annot_kws={"size": 12},fmt="d",ax=axTe)# font size, format in
digit

#Suggestion 4.Label confusion matrix heatmap with actual and predicted labels.
axTe.set_xticklabels(labels)
axTe.set_yticklabels(labels)
pltTe.title("Test confusion matrix")
pltTe.xlabel("Predicted")
pltTe.ylabel("Actual")
pltTe.show()

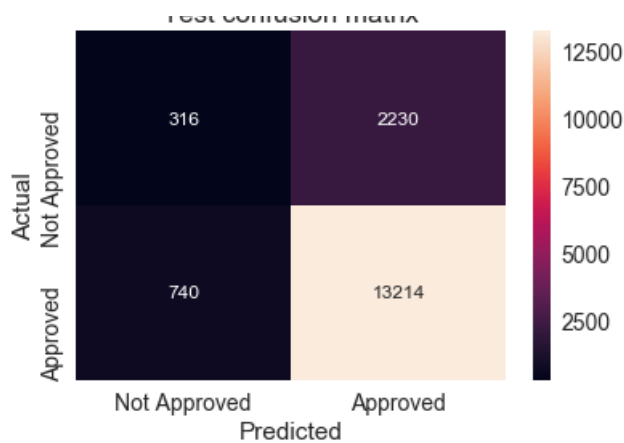
```

the maximum value of  $tpr \cdot (1 - fpr)$  0.11135206899920401 for threshold 0.855



the maximum value of  $tpr \cdot (1 - fpr)$  0.1175342031122122 for threshold 0.855

Test confusion matrix



1. Once after you plot the confusion matrix with the test data, get all the 'false positive data points'

- Plot the WordCloud [WordCloud](#)
- Plot the box plot with the 'price' of these 'false positive data points'
- Plot the pdf with the 'teacher\_number\_of\_previously\_posted\_projects' of these 'false positive data points'

In [164]:

```
predict_tfidf_W2V = predict(y_test_tfidf_w2v_pred, te_thresholds, test_fpr, test_tpr) # <==

#converting predict_bow list into numpy array.
# https://likegeeks.com/numpy-array-tutorial/
import numpy as np
y_predict_tfidf_W2V = np.array(predict_tfidf_W2V)
#print(y_predict_tfidf_W2V)
#print(type(y_predict_tfidf_W2V))

# traversing y_test and y_predict, and created new list y_FP, which is 1 for FALSE_POSITIVE, and 0
for rest
y_tfidf_W2V_FP=[]
count_tfidf_W2V_FP=0
# https://www.geeksforgeeks.org/numpy-iterating-over-array/
for ac,pre in np.nditer([y_test,y_predict_tfidf_W2V]):
    if(ac==0 and pre==1):
        y_tfidf_W2V_FP.append(1)
        count_tfidf_W2V_FP+=1
    else:
        y_tfidf_W2V_FP.append(0)

print("False Positive count:",count_tfidf_W2V_FP)
##print(type(y_test))
##print(y_test)
##print(type(y_predict_tfidf_W2V))

#print(y_predict)
#print(type(y_FP))
#print(y_FP)
#converting predict_bow list into numpy array.
y_tfidf_W2V_FP=np.array(y_tfidf_W2V_FP)
print(type(y_tfidf_W2V_FP))
print(y_tfidf_W2V_FP)
```

the maximum value of tpr\*(1-fpr) 0.1175342031122122 for threshold 0.855

False Positive count: 2230

<class 'numpy.ndarray'>

[0 0 0 ... 1 1 0]

In [165]:

```
X_test_tfidf_W2V_working = X_test.copy(deep=True) #<==

#print("X_test_working length:",len(X_test_working))
#print(type(X_test))
#print(X_test.shape)
#print(type(X_test_tfidf_W2V_working))
#print(X_test_tfidf_W2V_working#.head(2))
```

```
# Adding 3 numpy array into dataframe
#print(type(y_test))
#print(type(y_predict_tfidf_W2V))
#print(type(y_tfidf_W2V_FP))
#print(y_test)
#print(y_predict_tfidf_W2V)
#print(y_tfidf_W2V_FP)
#
#print(type(X_test_tfidf_W2V_working))
X_test_tfidf_W2V_working['Y_Actual'] = y_test
print(X_test_tfidf_W2V_working.head(2))
```

```

      Unnamed: 0      id      teacher_id school_state \
33081      141937  p034157  0d7b3cd172c5b19f83a0ed303f46b729      AR
26184      33737  p225681  d44f8cb33de45f83a0ed303f46b729      TN

      project_submitted_datetime \
33081      2016-09-26 16:20:26
26184      2016-12-26 14:31:27

      project_title \
33081  Scratching the Surface in Collaborative Learning!
26184      Ready 2 Read!

      project_essay_1 \
33081  With the upcoming school year, I am fortunate ...
26184  I am truly blessed to have the opportunity to ...

      project_essay_2 project_essay_3 \
33081  As my students become more technologically inc...      NaN
26184  My classroom consists of students on all learn...      NaN

      project_essay_4      ...      compound      price      quantity      essay_wc      essay_len \
33081      NaN      ...      0.9944      438.89      2      213      1528
26184      NaN      ...      0.9865      232.49      6      107      763

      title_wc      title_len      prj_res_sum_wc      prj_res_sum_len      Y_Actual
33081      4      40      14      106      1
26184      3      12      9      61      1

[2 rows x 31 columns]
```

In [166]:

```
# how to add numpy array into dataframe | https://stackoverflow.com/questions/26666919/add-column-
in-dataframe-from-list
print(X_test_tfidf_W2V_working.columns)
X_test_tfidf_W2V_working['y_predict_tfidf_W2V'] = y_predict_tfidf_W2V
X_test_tfidf_W2V_working['y_tfidf_W2V_FP'] = y_tfidf_W2V_FP
print(X_test_tfidf_W2V_working.columns)
```

```
Index(['Unnamed: 0', 'id', 'teacher_id', 'school_state',
      'project_submitted_datetime', 'project_title', 'project_essay_1',
      'project_essay_2', 'project_essay_3', 'project_essay_4',
      'project_resource_summary',
      'teacher_number_of_previously_posted_projects', 'project_is_approved',
      'clean_categories', 'clean_subcategories', 'clean_grade',
      'clean_tea_pfx', 'essay', 'neg', 'pos', 'neu', 'compound', 'price',
      'quantity', 'essay_wc', 'essay_len', 'title_wc', 'title_len',
      'prj_res_sum_wc', 'prj_res_sum_len', 'Y_Actual'],
      dtype='object')
Index(['Unnamed: 0', 'id', 'teacher_id', 'school_state',
      'project_submitted_datetime', 'project_title', 'project_essay_1',
      'project_essay_2', 'project_essay_3', 'project_essay_4',
      'project_resource_summary',
      'teacher_number_of_previously_posted_projects', 'project_is_approved',
      'clean_categories', 'clean_subcategories', 'clean_grade',
      'clean_tea_pfx', 'essay', 'neg', 'pos', 'neu', 'compound', 'price',
      'quantity', 'essay_wc', 'essay_len', 'title_wc', 'title_len',
      'prj_res_sum_wc', 'prj_res_sum_len', 'Y_Actual', 'y_predict_tfidf_W2V',
      'y_tfidf_W2V_FP'],
      dtype='object')
```

In [167]:

```
X_test_tfidf_W2V_FP=X_test_tfidf_W2V_working.copy(deep=True)
X_test_tfidf_W2V_FP.teacher_number_of_previously_posted_projects
X_test_tfidf_W2V_FP.price
X_test_tfidf_W2V_FP.columns
```

Out[167]:

```
Index(['Unnamed: 0', 'id', 'teacher_id', 'school_state',
      'project_submitted_datetime', 'project_title', 'project_essay_1',
      'project_essay_2', 'project_essay_3', 'project_essay_4',
      'project_resource_summary',
      'teacher_number_of_previously_posted_projects', 'project_is_approved',
      'clean_categories', 'clean_subcategories', 'clean_grade',
      'clean_tea_pfx', 'essay', 'neg', 'pos', 'neu', 'compound', 'price',
      'quantity', 'essay_wc', 'essay_len', 'title_wc', 'title_len',
      'prj_res_sum_wc', 'prj_res_sum_len', 'Y_Actual', 'y_predict_tfidf_W2V',
      'y_tfidf_W2V_FP'],
      dtype='object')
```

In [168]:

```
X_test_tfidf_W2V_FP=X_test_tfidf_W2V_FP.loc[:,
["teacher_number_of_previously_posted_projects","essay","price","y_tfidf_W2V_FP"]]
print(X_test_tfidf_W2V_FP.head(2))

X_test_tfidf_W2V_FP = X_test_tfidf_W2V_FP[X_test_tfidf_W2V_FP.y_tfidf_W2V_FP==True]
```

```
teacher_number_of_previously_posted_projects \
33081                                         43
26184                                         1

essay price \
33081 With the upcoming school year, I am fortunate ... 438.89
26184 I am truly blessed to have the opportunity to ... 232.49

y_tfidf_W2V_FP
33081      0
26184      0
```

In [169]:

```
# https://www.geeksforgeeks.org/generating-word-cloud-python/
# Python program to generate WordCloud

# importing all necessary modules
from wordcloud import WordCloud, STOPWORDS
import matplotlib.pyplot as plt
import pandas as pd

# Reads 'Youtube04-Eminem.csv' file
df = pd.read_csv(r"Youtube04-Eminem.csv", encoding="latin-1")

comment_words = ' '
stopwords = set(STOPWORDS)

# iterate through the csv file
for val in X_test_tfidf_W2V_FP["essay"][:1]:

    # typecaste each val to string
    val = str(val)

    # split the value
    tokens = val.split()

    # Converts each token into lowercase
    for i in range(len(tokens)):
        tokens[i] = tokens[i].lower()

    for words in tokens:
        comment_words = comment_words + words + ' '

wordcloud = WordCloud(width = 800, height = 800,
```

```

        background_color='white',
        stopwords = stopwords,
        min_font_size = 10).generate(comment_words)

# plot the WordCloud image
plt.figure(figsize = (8, 8), facecolor = None)
plt.imshow(wordcloud)
plt.axis("off")
plt.tight_layout(pad = 0)

plt.show()

```

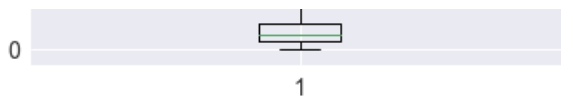


**Plot the box plot with the price of these false positive data points**

In [170]:

```
import matplotlib.pyplot as plt
# https://glowingpython.blogspot.com/2012/09/boxplot-with-matplotlib.html
plt.boxplot([X_test_tfidf_W2V_FP["price"]])
plt.title('Box Plots with the price of these tfidf_W2V\'s false positive data points')
#labels = ('Price')
#plt.xticks([1],labels,rotation=90)
plt.ylabel('Price')
plt.grid(True)
plt.show()
```





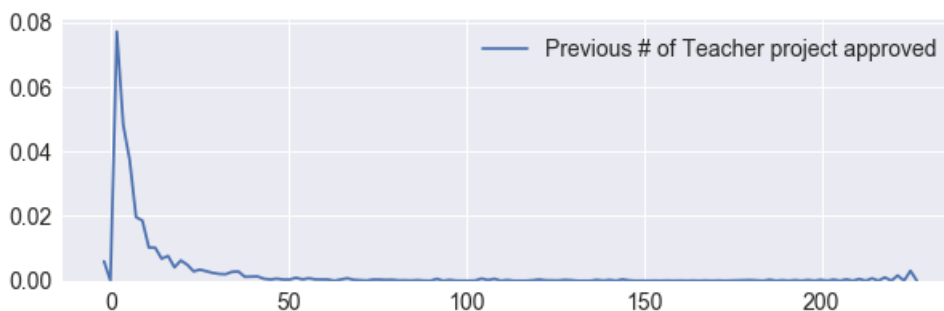
**Plot the pdf with the teacher\_number\_of\_previously\_posted\_projects of these false positive data points**

In [171]:

```
plt.figure(figsize=(10,3))

# https://seaborn.pydata.org/generated/seaborn.kdeplot.html | kernel density estimate | sns.kdeplot
# bw : {'scott' | 'silverman' | scalar | pair of scalars }, optional
# Name of reference method to determine kernel size, scalar factor, or scalar for each dimension
# of the bivariate plot. Note that the underlying computational libraries have different
# interpretations
# for this parameter: statsmodels uses it directly, but scipy treats it as a scaling factor
# for the standard deviation of the data.

sns.kdeplot(X_test_tfidf_W2V_FP["teacher_number_of_previously_posted_projects"],label="Previous #
of Teacher project approved", bw=0.6)
plt.legend()
plt.show()
```



## 2.5 [Task-2]Getting top 5k features using `feature\_importances\_`

In [172]:

```
# please write all the code with proper documentation, and proper titles for each subsection
# go through documentations and blogs before you start coding
# first figure out what to do, and then think about how to do.
# reading and understanding error messages will be very much helpfull in debugging your code
# when you plot any graph make sure you use
# a. Title, that describes your plot, this will be very helpful to the reader
# b. Legends if needed
# c. X-axis label
# d. Y-axis label
```

### 1. [Task-2]

- Select 5k best features from features of **Set 2** using `feature_importances_`, discard all the other remaining features and then apply any of the model of you choice i.e. (Decision tree, Logistic Regression, Linear SVM), you need to do hyperparameter tuning corresponding to the model you selected and procedure in step 2 and step 3

In [217]:

```
# merge two sparse matrices: https://stackoverflow.com/a/19710648/4084039
from scipy.sparse import hstack
X_tr_tfidf = hstack((X_train_text_tfidf, X_train_title_tfidf, X_train_state_ohe, X_train_clean_ohe,
X_train_cleanSub_ohe, X_train_grade_ohe, X_train_teacher_ohe, X_train_prjResSum_ohe,
X_train_quantity_norm, X_train_TprevPrj_norm, X_train_price_norm)).tocsr()
X_te_tfidf = hstack((X_test_text_tfidf, X_test_title_tfidf, X_test_state_ohe, X_test_clean_ohe, X_
test_cleanSub_ohe, X_test_grade_ohe, X_test_teacher_ohe, X_test_prjResSum_ohe,
X_test_quantity_norm, X_test_TprevPrj_norm, X_test_price_norm)).tocsr()
```

```

print("Final Data matrix | tfidf")
print(X_tr_tfidf.shape, y_train.shape)
print(X_te_tfidf.shape, y_test.shape)
print("="*100)

```

```

Final Data matrix | tfidf
(33500, 19119) (33500,)
(16500, 19119) (16500,)
=====

```

In [218]:

```

modelT2 = DecisionTreeClassifier(max_depth=50)
modelT2.fit(X_tr_tfidf, y_train)

```

Out[218]:

```

DecisionTreeClassifier(class_weight=None, criterion='gini', max_depth=50,
                        max_features=None, max_leaf_nodes=None,
                        min_impurity_decrease=0.0, min_impurity_split=None,
                        min_samples_leaf=1, min_samples_split=2,
                        min_weight_fraction_leaf=0.0, presort=False, random_state=None,
                        splitter='best')

```

In [219]:

```

print(type(modelT2.feature_importances_))
print(len(modelT2.feature_importances_))
print(X_tr_tfidf.shape)
Impfeature = modelT2.feature_importances_
#Impfeature[13001:14000]
print(len(Impfeature))

```

```

import numpy as npcont
count=0
for i in np.nditer(Impfeature):
    if i != 0:
        count+=1
        #print(count)
        #print(i)
print("Total:",count)

```

```

<class 'numpy.ndarray'>
19119
(33500, 19119)
19119
Total: 1048

```

In [220]:

```

#print(type(X_tr_tfidf))
#print(X_tr_tfidf)
print(type(Impfeature.argsort()))
print(Impfeature.argsort())

ImpFeatureSorted = Impfeature.argsort()
ImpFeatureSorted=ImpFeatureSorted[::-1]
print(ImpFeatureSorted.shape)
print(ImpFeatureSorted)

```

```

<class 'numpy.ndarray'>
[  0 12710 12709 ... 4231 13842 16981]
(19119,)
[16981 13842 4231 ... 12709 12710    0]

```

In [221]:

```

X_tr_tfidf=X_tr_tfidf[:,ImpFeatureSorted[:5000]]

```



In [222]:

```
X_tr_tfidf.shape
```

Out[222]:

```
(33500, 5000)
```

In [223]:

```
X_te_tfidf=X_te_tfidf[:,ImpFeatureSorted[:5000]]
X_te_tfidf.shape
```

Out[223]:

```
(16500, 5000)
```

## Logistic Regressions

In [224]:

```
#code source:
http://occam.olin.edu/sites/default/files/DataScienceMaterials/machine_learning_lecture_2/Machine%20Learning%20Lecture%202.html
from sklearn.model_selection import train_test_split
#from sklearn.grid_search import GridSearchCV
from sklearn.model_selection import GridSearchCV
from sklearn.datasets import *
from sklearn.linear_model import LogisticRegression
```

```
c_range=[10**-5, 10**-4, 10**-2, 10**0, 10**2, 10**4, 10**5]
param_grid=dict(C=c_range)
```

```
#Using GridSearchCV
# https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.GridSearchCV.html
modelLR = GridSearchCV(LogisticRegression(class_weight="balanced"), param_grid, scoring = 'f1', cv
=5)
modelLR.fit(X_tr_tfidf, y_train)

print(modelLR.best_estimator_)
print(modelLR.score(X_te_tfidf, y_test))
```

```
LogisticRegression(C=1, class_weight='balanced', dual=False,
                    fit_intercept=True, intercept_scaling=1, max_iter=100,
                    multi_class='ovr', n_jobs=1, penalty='l2', random_state=None,
                    solver='liblinear', tol=0.0001, verbose=0, warm_start=False)
0.7948330683624801
```

In [225]:

```
train_tf_auc= modelLR.cv_results_['mean_train_score']
train_tf_auc_std= modelLR.cv_results_['std_train_score']
cv_tf_auc = modelLR.cv_results_['mean_test_score']
cv_tf_auc_std= modelLR.cv_results_['std_test_score']

CC = []
from math import log
CC = [np.log10(x) for x in c_range]
print(CC)

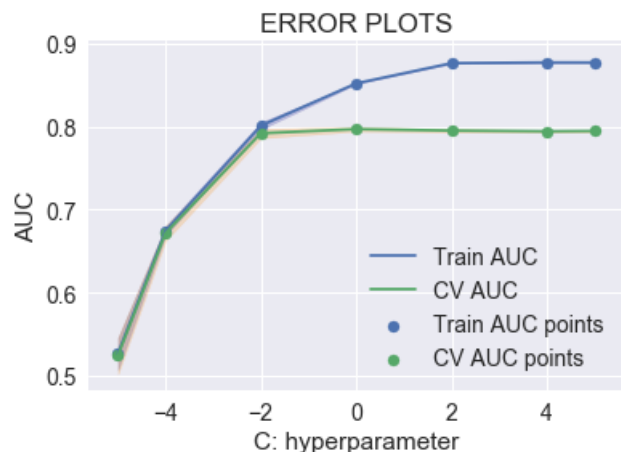
plt.plot(CC, train_tf_auc, label='Train AUC')
# this code is copied from here: https://stackoverflow.com/a/48803361/4084039
plt.gca().fill_between(CC,train_tf_auc - train_tf_auc_std,train_tf_auc + train_tf_auc_std,alpha=0.2
,color='darkblue')

plt.plot(CC, cv_tf_auc, label='CV AUC')
# this code is copied from here: https://stackoverflow.com/a/48803361/4084039
plt.gca().fill_between(CC,cv_tf_auc - cv_tf_auc_std,cv_tf_auc + cv_tf_auc_std,alpha=0.2,color='dark
orange')
```

```
plt.scatter(CC, train_tf_auc, label='Train AUC points')
plt.scatter(CC, cv_tf_auc, label='CV AUC points')
```

```
plt.legend()
plt.xlabel("C: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid(True)
plt.show()
```

```
[-5.0, -4.0, -2.0, 0.0, 2.0, 4.0, 5.0]
```



In [226]:

```
best_tuned_parameters = [{'C': [1]}]
```

In [227]:

```
# https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_curve.html#sklearn.metrics.roc_curve
from sklearn.metrics import roc_curve, auc

modelLRb = GridSearchCV(LogisticRegression(), best_tuned_parameters)
modelLRb.fit(X_tr_tfidsf, y_train)
# roc_auc_score(y_true, y_score) the 2nd parameter should be probability estimates of the positive class
# not the predicted outputs

print(modelLRb.best_estimator_)
print(modelLRb.score(X_te_tfidsf, y_test))

y_train_tf_pred = modelLRb.predict_proba(X_tr_tfidsf)[:,1]
y_test_tf_pred = modelLRb.predict_proba(X_te_tfidsf)[:,1]

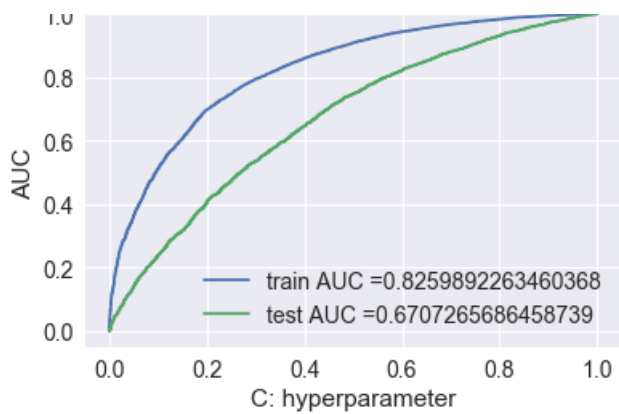
train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, y_train_tf_pred)
test_fpr, test_tpr, te_thresholds = roc_curve(y_test, y_test_tf_pred)

plt.plot(train_fpr, train_tpr, label="train AUC =" + str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="test AUC =" + str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel("C: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.grid(True)
plt.show()
```

```
LogisticRegression(C=1, class_weight=None, dual=False, fit_intercept=True,
                    intercept_scaling=1, max_iter=100, multi_class='ovr', n_jobs=1,
                    penalty='l2', random_state=None, solver='liblinear', tol=0.0001,
                    verbose=0, warm_start=False)
0.8407272727272728
```

ERROR PLOTS





In [228]:

```
# https://scikit-learn.org/stable/modules/generated/sklearn.metrics.confusion_matrix.html
print("="*100)
from sklearn.metrics import confusion_matrix
print("Train confusion matrix")
print(confusion_matrix(y_train, predict(y_train_tf_pred, tr_thresholds, train_fpr, train_tpr)))
print("Test confusion matrix")
print(confusion_matrix(y_test, predict(y_test_tf_pred, te_thresholds, test_fpr, test_tpr)))
```

```
=====

Train confusion matrix
the maximum value of tpr*(1-fpr) 0.5630193014364399 for threshold 0.829
[[ 3835  1333]
 [ 6836 21496]]
Test confusion matrix
the maximum value of tpr*(1-fpr) 0.38954756065857055 for threshold 0.854
[[1513 1033]
 [4807 9147]]
```

In [229]:

```
import seaborn as snTr
import seaborn as snTe
import pandas as pdH
import matplotlib.pyplot as pltTr
import matplotlib.pyplot as pltTe

# https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
arrayTr=confusion_matrix(y_train, predict(y_train_tf_pred, tr_thresholds, train_fpr, train_tpr))
df_cmTr = pdH.DataFrame(arrayTr,range(2),range(2))
#print(arrayTr)
# https://stackoverflow.com/questions/32723798/how-do-i-add-a-title-to-seaborn-heatmap
axTr = pltTr.axes()

snTr.set(font_scale=1.4)#for label size

# https://seaborn.pydata.org/generated/seaborn.heatmap.html

snTr.heatmap(df_cmTr, annot=True,annot_kws={"size": 12},fmt="d",ax=axTr)# font size, format in
digit

labels=['Not Approved','Approved']
axTr.set_xticklabels(labels)
axTr.set_yticklabels(labels)
#Suggestion 4.Label confusion matrix heatmap with actual and predicted labels.
pltTr.title("Train confusion matrix")
pltTr.xlabel("Predicted")
pltTr.ylabel("Actual")
pltTr.show()

# https://stackoverflow.com/questions/50947776/plot-two-seaborn-heatmap-graphs-side-by-side
#fig, ax =plt.subplots(1,1)

# https://stackoverflow.com/questions/35572000/how-can-i-plot-a-confusion-matrix
arrayTe=confusion_matrix(y_test, predict(y_test_tf_pred, te_thresholds, test_fpr, test_tpr))
df_cmTe = pdH.DataFrame(arrayTe,range(2),range(2))
```

```

axTe = pltTe.axes()

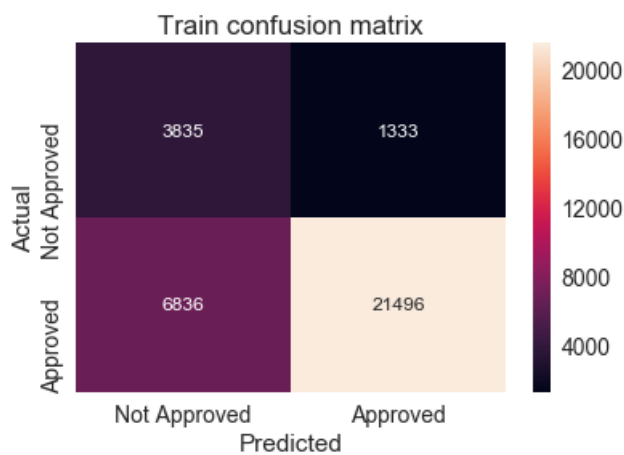
snTe.set(font_scale=1.4)#for label size

# https://seaborn.pydata.org/generated/seaborn.heatmap.html
snTe.heatmap(df_cmTe, annot=True,annot_kws={"size": 12},fmt="d",ax=axTe)# font size, format in
digit

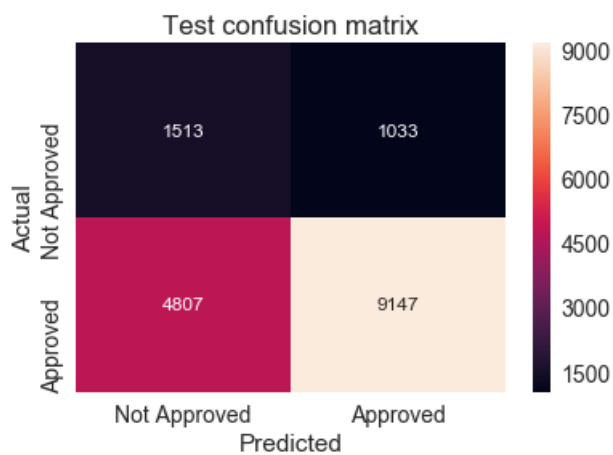
#Suggestion 4.Label confusion matrix heatmap with actual and predicted labels.
axTe.set_xticklabels(labels)
axTe.set_yticklabels(labels)
pltTe.title("Test confusion matrix")
pltTe.xlabel("Predicted")
pltTe.ylabel("Actual")
pltTe.show()

```

the maximum value of  $tpr*(1-fpr)$  0.5630193014364399 for threshold 0.829



the maximum value of  $tpr*(1-fpr)$  0.38954756065857055 for threshold 0.854



### 3. Conclusion

In [231]:

```

# Please compare all your models using Prettytable library
from prettytable import PrettyTable

x = PrettyTable()

x.field_names = ["Vectorizer", "Algorithm", "Max Depth", "min_smaple_split", "Test AUC", "Best score"]

x.add_row(["BOW", "Decision Tree", 1, 5, 0.53554243,0.8976886471787899 ])
x.add_row(["TFIDF", "Decision Tree", 500, 5, 0.52738914,0.7492121212121212 ])
x.add_row(["AVG W2V", "Decision Tree", 1, 5, 0.5355424359,0.8456969696969697 ])
x.add_row(["TFIDF W2V", "Decision Tree", 1, 5, 0.5355424359,0.8456969696969697 ])

```

```
x.add_row(["TFIDF", "Logistic Regression", "C=1", " ", 0.9707265686, 0.84072727272728 ])

print(x)
```

Vectorizer	Algorithm	Max Depth	min_sample_split	Test AUC	Best score
BOW	Decision Tree	1	5	0.53554243	0.8976886471787899
TFIDF	Decision Tree	500	5	0.52738914	0.7492121212121212
AVG W2V	Decision Tree	1	5	0.5355424359	0.8456969696969697
TFIDF W2V	Decision Tree	1	5	0.5355424359	0.8456969696969697
TFIDF	Logistic Regression	C=1		0.9707265686	0.8407272727272728

## Summary

### Step followed

- Preprocessing of Project\_subject\_categories Project\_subject\_subcategories project\_grade\_category teacher\_prefix Project\_essay Project\_title project\_resource\_summary
- Numeric feature for Text no of words in essay lenght of each cell in essay no of words in Title lenght of each cell in Title no of words in Project resource summary lenght of each cell in Project resource summary
- Using Pretrained Models: Avg W2V
- Computing Sentiment Scores for Project essay. Added below columns neg pos neu compound
- Added all the features to project\_data
- Took data points for doing the assignment and separate the Class lable (Project\_is\_approved)
- Splitting Data into Train and Test.
- Making datamodel ready

##### text

- encoding of school\_state is splited into Train and Test vector and stored the feature name in aa
- encoding of clean\_category is splited into Train and Test vector and stored the feature name in b
- encoding of clean\_subcategory is splited into Train and Test vector and stored the feature name in c
- encoding of project\_grade\_category is splited into Train and Test vector and stored the feature name in d
- encoding of teacher\_prefix is splited into Train and Test vector and stored the feature name in e
- encoding of project\_resource\_summary is splited into Train and Test vector

r

and stored the feature name in ff

##### numeric

- encoding of quantity is splited into Train and Test vector
- encoding of teacher\_number\_of\_previously\_posted\_projects is splited into Train and Test vector
- encoding of price is splited into Train and Test vector and stored (quantity, teacher\_number\_of\_previously\_posted\_projects, price) the feature name in h
- encoding of sentimental score | neg, is splited into Train and Test vector

r

```

r      - encoding of sentimental score | pos, is splited into Train and Test vecto
r      - encoding of sentimental score | neu, is splited into Train and Test vecto
r      - encoding of sentimental score | compound, is splited into Train and Test
vector
      - encoding of numerical | number of words in the title, is splited into Tra
in and Test vector
      - encoding of numerical | number of words in the essay, is splited into Tra
in and Test vector
      - encoding of project_essay(BOW) is splited into Train and Test vector
        and stored the feature name in g
      - encoding of project_title(BOW) is splited into Train and Test vector
        and stored the feature name in k
      - encoding of project_essay(TFIDF) is splited into Train and Test vector
        and stored the feature name in ii
      - encoding of project_title(TFIDF) is splited into Train and Test vector
        and stored the feature name in j
      - encoding of project_essay(AVG W2V) is splited into Train and Test vector
      - encoding of project_title(AVG W2V) is splited into Train and Test vector
      - encoding of project_essay(TFIDF W2V) is splited into Train and Test vecto
r
r      - encoding of project_title(TFIDF W2V) is splited into Train and Test vecto
r

```

## For SET 1

### Merging all the above features for SET 1

- Horizontally merging( with hstack) all categorical, numerical features + project\_title(BOW) + preprocessed\_essay (BOW)
- Fit a model on on train (on above merge features) data by using  
GridSearchCV(DecisionTreeClassifier(class\_weight="balanced"))
- Draw a graph in Train and CV for varies values of alpha
- Take Best\_Alpha by *bestestimator* and draw graph for Test\_AUC
- Create Confusion matrix, in heatmap.
- create a wordCloud
  - created a list of False positive, y\_FP
  - deep copied x\_test into x\_test\_working
  - Appended x\_test\_working with, Y\_Actual, Y\_Predict, Y\_FP
  - deep copied x\_test\_working into x\_test\_FP
  - New x\_test\_FP with only column as [teacher\_number\_of\_previously\_posted\_projects", "essay", "price", "Y\_FP]
  - Selected only those rows, which have FP coulmn as True.
  - created Wordcloud graph
  - Created the box plot with the price and False positive points
  - Plot PDF of Teacher project approved with False positive points.

### Graphviz

- added all the feature\_list(aa,b,c,d,e,ff,g,k and h)
- run export\_graphviz()

## For SET 2

### Merging all the above features for SET 2

- Horizontally merging( with hstack) all categorical, numerical features + project\_title(TFIDF) + preprocessed\_essay (TFIDF)
- Fit a model on on train (on above merge features) data by using  
ridSearchCV(DecisionTreeClassifier(class\_weight="balanced"))
- Draw a graph in Train and CV for varies values of alpha
- Take Best\_Alpha by *bestestimator* and draw graph for Test\_AUC
- Create Confusion matrix, in heatmap.
- create a wordCloud
  - created a list of False positive, y\_tfidf\_FP
  - deep copied x\_test into X\_test\_tfidf\_working
  - Appended X\_test\_tfidf\_working with, Y\_Actual, y\_predict\_tfidf, y\_tfidf\_FP
  - deep copied X\_test\_tfidf\_working into x\_test\_FP

- New x\_test\_FP with only column as [teacher\_number\_of\_previously\_posted\_projects","essay","price","y\_tfidf\_FP]
- Selected only those rows, which have FP columns as True.
- created Wordcloud graph
- Created the box plot with the price and False positive points
- Plot PDF of Teacher project approved with False positive points.

#### Graphviz

- added all the feature\_list(aa,b,c,d,e,ff,ii,j and h)
- run export\_graphviz()

## For SET 3

### Merging all the above features for SET 3

- Horizontally merging( with hstack) all categorical, numerical features + project\_title(AVG W2V)) + preprocessed\_essay (AVG W2V))
- Fit a model on on train (on above merge features) data by using GridSearchCV(DecisionTreeClassifier(class\_weight="balanced"))
- Draw a graph in Train and CV for varies values of alpha
- Take Best\_Alpha by *bestestimator* and draw graph for Test\_AUC
- Create Confusion matrix, in heatmap.
- create a wordCloud
  - created a list of False positive, y\_avgW2V\_FP
  - deep copied x\_test into X\_test\_avgW2V\_working
  - Appended X\_test\_avgW2V\_working with, Y\_Actual, y\_predict\_avgW2V, y\_avgW2V\_FP
  - deep copied X\_test\_avgW2V\_working into x\_test\_FP
  - New x\_test\_FP with only column as [teacher\_number\_of\_previously\_posted\_projects","essay","price","y\_avgW2V\_FP]
  - Selected only those rows, which have FP columns as True.
  - created Wordcloud graph
  - Created the box plot with the price and False positive points
  - Plot PDF of Teacher project approved with False positive points.

## For SET 4

### Merging all the above features for SET 4

- Horizontally merging( with hstack) all categorical, numerical features + project\_title(TFIDF W2V)) + preprocessed\_essay (TFIDF W2V))
- Fit a model on on train (on above merge features) data by using GridSearchCV(DecisionTreeClassifier(class\_weight="balanced"))
- Draw a graph in Train and CV for varies values of alpha
- Take Best\_Alpha by *bestestimator* and draw graph for Test\_AUC
- Create Confusion matrix, in heatmap.
- create a wordCloud
  - created a list of False positive, y\_tfidf\_W2V\_FP
  - deep copied x\_test into X\_test\_tfidf\_W2V\_working
  - Appended X\_test\_tfidf\_W2V\_working with, Y\_Actual, y\_predict\_tfidf\_W2V, y\_tfidf\_W2V\_FP
  - deep copied X\_test\_tfidf\_W2V\_working into X\_test\_tfidf\_W2V\_FP
  - New X\_test\_tfidf\_W2V\_FP with only column as [teacher\_number\_of\_previously\_posted\_projects","essay","price","y\_tfidf\_W2V\_FP]
  - Selected only those rows, which have FP columns as True.
  - created Wordcloud graph
  - Created the box plot with the price and False positive points
  - Plot PDF of Teacher project approved with False positive points.

## Task 2, for selecting featureimportantes for 5K

- fit the model in DecisionClassifier(max\_depth=50)
- take the top 5000 important\_features
- Take the indexes of importatn features, by argsort() command, and sort it.
- for sorted index, we get in above steps, takes corresponding all rows from TFIDF's train and test data
- Run a logistice regression on it
- Draw a graph in Train and CV for varies values of alpha
- Take Best\_Alpha by *bestestimator* and draw graph for Test\_AUC
- Create Confusion matrix, in heatmap.

- Create confusion matrix, in heatmap.